

1 **For biorxiv**

2

3 **16 Oct, 2019.**

4

5 *Deaminase associated single nucleotide variants in blood and*  
6 *saliva-derived exomes from healthy subjects*

7

8

9 Nathan E. Hall, Jared Mamrot, Christopher M.A. Frampton, Prue Read, Edward J. Steele, Robert J.  
10 Bischof, and Robyn A. Lindley

11

12

13 **Author Information:**

14

15 Nathan E. Hall                      nathan.hall@gmdxgroup.com

16 Jared Mamrot                        jared.mamrot@gmdxgroup.com

17 Christopher M.A. Frampton      chris.frampton@otago.ac.nz

18 Prue Read                            pruef@optusnet.com.au

19 Edward J. Steele                    e.j.steele@bigpond.com

20 Robert J. Bischof                    robert.bischof@hudson.org.au

21 Robyn A. Lindley                    robyn.lindley@gmdxgroup.com

22

23

24 *Deaminase associated single nucleotide variants in blood and*  
25 *saliva-derived exomes from healthy subjects*

26

27 Nathan E. Hall<sup>1,2,\*</sup>, Jared Mamrot<sup>1,3,4</sup>, Christopher M.A. Frampton<sup>5</sup>, Prue Read<sup>1,6</sup>, Edward J. Steele<sup>7,8</sup>,  
28 Robert J. Bischof<sup>3</sup>, and Robyn A. Lindley<sup>1,9</sup>

29

30 <sup>1</sup>GMDx Group Ltd, Melbourne, 3000, Victoria, AUSTRALIA;

31 <sup>2</sup>Department of Animal, Plant and Soil Sciences, La Trobe University, Melbourne, Victoria,  
32 AUSTRALIA;

33 <sup>3</sup>The Ritchie Centre, Hudson Institute of Medical Research, Clayton, 3168 Victoria, AUSTRALIA;

34 <sup>4</sup>Monash University, Clayton, Victoria, AUSTRALIA;

35 <sup>5</sup>Department of Medicine, University of Otago, Christchurch, NEW ZEALAND;

36 <sup>6</sup>Five Corners Pty Ltd 13/76 Reserve Road, Artarmon, NSW AUSTRALIA;

37 <sup>7</sup>CYO'Connor ERADE Village Foundation, 24 Genomics Rise, Piara Waters, Perth, AUSTRALIA;

38 <sup>8</sup>Melville Analytics Pty Ltd, Melbourne, Victoria, AUSTRALIA;

39 <sup>9</sup>Department of Clinical Pathology, The Victorian Comprehensive Cancer Centre, Faculty of Medicine,  
40 Dentistry & Health Sciences, University of Melbourne, Victoria, AUSTRALIA.

41

42 **\*Corresponding Author** Dr. Nathan E. Hall, GMDx Group Ltd, 162 Collins Street, Melbourne Vic,  
43 3000, AUSTRALIA email: [nathan.hall@gmdxgroup.com](mailto:nathan.hall@gmdxgroup.com)

44

45 **Running head:** *Deaminase SNVs in Healthy Subjects*

46

47 **Key words:** Blood-saliva concordance, Cytosine and Adenosine Deamination, Somatic Mutation,  
48 Single Nucleotide Variations, Whole Exome Sequencing, AID, APOBEC and ADAR Deamination

49

50

51 **Abstract**

52 **Background:**

53 Deaminases play an important role in shaping inherited and somatic variants. Disease related SNVs are  
54 associated with deaminase mutagenesis and genome instability. Here, we investigate the reproducibility  
55 and variance of whole exome SNV calls in blood and saliva of healthy subjects and analyze variants  
56 associated with AID, ADAR, APOBEC3G and APOBEC3B deaminase sequence motifs.

57 **Methods:**

58 Samples from twenty-four healthy Caucasian volunteers, allocated into two groups, underwent whole  
59 exome sequencing. Group 1 (n=12) analysis involved one blood and four saliva replicates. A single  
60 saliva sample was sequenced for Group 2 subjects (n=12). Overall, a total of 72 whole exome datasets  
61 were analyzed. Biological (Group 1 & 2) and technical (Group 1) variance of SNV calls and deaminase  
62 metrics were calculated and analyzed using intraclass correlation coefficients. Candidate somatic SNVs  
63 were identified and evaluated.

64 **Results:**

65 We report high blood-saliva concordance in germline SNVs from whole exome sequencing.  
66 Concordant SNVs, found in all subject replicates, accounted for 97% of SNVs located within the  
67 protein coding sequence of genes. Discordant SNVs have a 30% overlap with variants that fail  
68 gnomAD quality filters and are less likely to be found in dbSNP. SNV calls and deaminase-associated  
69 metrics were found to be reproducible and robust (intraclass correlation coefficients >0.95). No somatic  
70 SNVs were conclusively identified when comparing blood and saliva samples.

71 **Conclusions:**

72 Saliva and blood both provide high quality sources of DNA for whole exome sequencing, with no  
73 difference in ability to resolve SNVs and deaminase-associated metrics. We did not identify somatic  
74 SNVs when comparing blood and saliva of healthy individuals, and we conclude that more specialized  
75 investigative methods are required to comprehensively assess the impact of deaminase activity on  
76 genome stability in healthy individuals.

77

## 78 **Background**

79 APOBEC/AID deaminases are a recognized endogenous source of genome instability [1–5]. Somatic  
80 mutations caused by deamination events have been identified in cancer *in vitro* and *in vivo* [6–9], and  
81 evidence of deaminase-associated mutations in non-cancerous conditions is emerging, such as various  
82 viral infections and neurodegenerative diseases [10,11]. Deaminases have also recently been implicated  
83 in accumulation of pre-cancerous mutations [12], and as a causative driver of many human SNPs [13].

84

85 Deaminases predominantly drive C-to-U(T) and A-to-I(G) transition mutations, however DNA repair  
86 mechanisms typically prevent deamination from compromising genome integrity and causing somatic  
87 mutation [14,15]. Pathophysiological processes can disrupt normal DNA repair, resulting in mosaic  
88 manifestation of deaminase-associated single nucleotide variants (SNVs) in affected tissues [16].

89 Although deaminases employ similar biochemical mechanisms, each has a unique binding domain  
90 associated with one or more DNA motifs [17,18]. Deaminase motifs can be identified and quantified in  
91 Next-Generation Sequencing (NGS) data facilitating diagnosis of the specific cause of the mutation.

92 For example, AID targets C-sites in the context of WRC motifs (W = A or T; R = A or G; reverse  
93 complements as GYW, with Y = T or C), APOBEC3G deaminates CC sites (or GG) and APOBEC3B  
94 deaminates TCW (or WGA) motifs and ADARs deaminate WA sites [2,19,20]. Establishing  
95 reproducible and robust deaminase-associated SNV profiles in healthy people will improve the utility

96 of mutation profiling techniques for monitoring progression of diseases such as cancer, and for  
97 understanding patient response to treatment.

98

99 Sampling of saliva or buccal cells is a widely employed technique for collecting human DNA for  
100 ancestry, forensic, medical and research purposes [21,22,23,24]. DNA extracted from saliva can be  
101 analyzed using various NGS techniques, however the quality of DNA derived from saliva can be  
102 compromised by metagenomic DNA and activity of various enzymes and antibacterial factors. There  
103 are several practical advantages to this DNA source, such as ease of sampling and additional  
104 sequencing information about metagenomic populations [25,26,27], however DNA obtained from  
105 saliva is not yet routinely used for detecting SNVs.

106

107 Here, we report profiles for SNVs associated with deaminase motifs for a cohort of 24 healthy human  
108 subjects using whole exome sequencing (WES). For twelve of these subjects (Group 1) we compare  
109 blood with biological and technical saliva replicates from Caucasian volunteers of different age groups  
110 and sex and hypothesize that deaminase-associated SNV profiles of a cohort of healthy individuals will  
111 show a high concordance between saliva and whole blood DNA in a reproducible and robust manner.

112

## 113 **Methods**

### 114 **Healthy subject selection**

115 In total, 24 healthy Caucasian subjects were recruited for this study. Volunteers were considered  
116 healthy if they had blood pressure and heart rate within normal ranges, had never smoked, were only  
117 light drinkers (<14 units of alcohol weekly), had no major viral infections or immune related diseases  
118 and did not take any regular medication. Eight subjects were recruited into each of the three age groups  
119 18-19, 30-39, and 50-59, with an equal ratio of males to females in each group. These subjects were

120 randomly allocated into two groups of equal sex and age group. Group 1 (n=12) involved analysis of  
121 blood and saliva sample replicates. Group 2 (n=12) involved analysis of saliva-1 sample only. This  
122 project was approved by the Monash Health Human Research Ethics Committee (16281L: “A study to  
123 measure the Targeted Somatic Mutation (TSM) test platform performance characteristics and evaluate  
124 its suitability for clinical use”).

125

### 126 **Sample collection**

127 For each subject, two saliva samples were collected, 30 minutes apart, using the Oragene DNA (OG-  
128 500) saliva collection kit. Whole blood samples were collected into sterile EDTA tubes.

129

### 130 **DNA extraction from saliva and whole blood**

131 DNA was extracted from samples using the QIAasympyony and the Qiagen DSP DNA Mini Kit. The  
132 extracted DNA was eluted in 100uL of Qiagen ATE buffer.

133

### 134 **Library preparation**

135 Whole exome sequencing library preparation was performed at the Monash Health Translation Precinct  
136 (MHTP) Medical Genomics Facility using the Agilent SureSelectXT Target Enrichment System  
137 according to protocol G7530- 90000, Version C0, December 2016. Capture Probes: Agilent SureSelect  
138 Clinical Research Exome Cat No 5190-7344; Design ID S06588914. Libraries were QC-checked using  
139 the Agilent BioAnalyzer and quantified with Qubit.

140

### 141 **Whole Exome Sequencing**

142 Four samples per lane were clustered on the c-bot using 200pM of library pool using Illumina Protocol  
143 15006165 v02 Jan 2016. Raw data was generated on the Illumina HiSeq 3000 with 100 base-pair  
144 paired-end (PE) sequencing with Illumina Protocol 15066493 Rev A, February 2015. Total PE reads

145 per sample were between 110 million and 183 million per exome, excluding HP\_4 saliva-1 which had  
146 82 million. The median number of PE reads per sample (137 million) and additional summary statistics  
147 are provided in Supplementary Table 1.

148

### 149 **Bioinformatics analysis**

150 WES read quality was assessed using FastQC (v0.11.7) [28]. Adapters were trimmed with cutadapt  
151 (v1.16) [29] and mapped to the human genome version hg19 with bwa (v0.7.13-r1126) [30] with the  
152 parameters “bwa mem -M -t 5 -k 19”. Duplicates were marked with Picard MarkDuplicates (v2.6.0)  
153 (<http://broadinstitute.github.io/picard>). Single Nucleotide Variant (SNV) calls were made with Strelka2  
154 (v2.8.4) [31] with default parameters using “configureStrelkaGermlineWorkflow.py -exome”. Variants  
155 failed quality filtering if they had a ConservativeGenotypeQuality < 15, a RelativeTotalLocusDepth <  
156 3, or a SampleStrandBias > 10. Variants remaining after quality filtering were converted from hg19 to  
157 hg38 coordinates using UCSC’s LiftOver tool [32]. Variants were identified as being located in the  
158 coding sequence (CDS) of genes according to Ensembl version 92 [33]. Candidate SNVs were com-  
159 pared against dbSNP v150 (10-07-2017) [34] and gnomAD exome release v2.0.2 [35]. Candidate so-  
160 matic SNVs were identified using Strelka2 with default parameters in the “configureStrelkaSomatic-  
161 Workflow.py -exome” pipeline. Unmapped reads were QC checked using FastQC and MultiQC (v1.6)  
162 [36] (<https://jpmam1.github.io/MultiQC>). To determine the source of the unmapped reads from repre-  
163 sentative subjects HP\_1 and HP\_2, these were aligned to the NCBI “non-redundant” (nr) database  
164 comprised of 4,348,972 protein sequences from eukaryotic and prokaryotic organisms, using DIA-  
165 MOND BLASTx (v0.9.22.123) [37]. Alignments were visualized using MEGAN6 (v6.12.2) [38,39].

166

### 167 **Deaminase motifs in WES data**

168 SNVs occurring within four key deaminase motifs were identified and quantified. The motifs used  
169 were AID: WRC / GYW, ADAR: WA / TW, APOBEC3G: CC / GG, APOBEC3B: TCW / WGA

170 where W=A or T, Y=C or T, R=A or G [2,19,20]. The base mutated in each motif is underlined. Motif  
171 searches were conducted according to the direction of the gene. Transition/transversion ratios (Ti/Tv)  
172 were calculated as the proportion of total transition variants. Strand bias was calculated as the  
173 proportion of variants on the forward strand (e.g. C:G and A:T as percentages). Motif-independent  
174 metrics and SNVs not associated with motifs of AID, ADAR, APOBEC3G or APOBEC3B (denoted  
175 “Other”) were also quantified.

176

### 177 **Experimental design**

178 Five WES datasets were generated for twelve subjects (Group 1), comprised of two males and two  
179 females from three age categories (18-19, 30-39 and 50-59). As described in Figure 1, replicates were  
180 generated from two saliva samples at the DNA extraction stage (saliva-1C) and at the library  
181 preparation stage (saliva-2A). These technical and biological replicates enabled analysis of  
182 concordance between replicates and provided a measure of technical variance and noise. This study  
183 design allows quantitative comparisons between blood and saliva, between saliva sample replicates and  
184 between technical saliva replicates at the DNA extraction and library preparation level. Group 2  
185 subjects (n=12) underwent WES of saliva-1 samples only and were used in the calculation of biological  
186 variance between subjects.

187 Mapped and unmapped WES reads were analyzed for genomic variants and off-target metagenomic  
188 contamination.

189

### 190 **Statistical analyses**

191 Intraclass correlation was calculated for SNV counts and deaminase motif metrics using the formula  
192 described by Shrout & Fleiss [40]:  $\frac{\sigma_{\alpha}^2}{\sigma_{\alpha}^2 + \sigma_{\epsilon}^2}$ . In brief,  $\sigma_{\alpha}$  represents the biological variance between  
193 subjects and  $\sigma_{\epsilon}$  represents the technical variance within subject replicates. Used here,  $\sigma_{\alpha}$  is the standard



194 deviation of saliva-1 samples across all 24 subjects (Group 1 and 2). For each Group 1 subject (n=12),  
195 a standard deviation is calculated and represents the variance within blood and saliva samples, DNA  
196 extraction and libraries;  $\sigma_e$  is the average of these twelve standard deviations.

197 SNVs termed *discordant* were found in 1, 2, 3, or 4 of the 5 samples, but not in all samples for each  
198 individual. Venn diagrams of concordant/discordant SNVs were generated at  
199 <http://bioinformatics.psb.ugent.be/webtools/Venn/>. Pairwise sample comparisons were conducted for  
200 discordant SNVs and analyzed using one-way ANOVA.

201

## 202 **Results**

### 203 **Blood and saliva whole exome sequencing**

204 Saliva and blood samples from 12 healthy volunteers, Group 1 subjects, underwent sequencing and  
205 analysis according to the workflow illustrated in Figure 1. In addition, 12 exomes were obtained from  
206 Saliva-1 samples from the remaining 12 recruited healthy volunteers, Group 2 subjects (Table 2). For  
207 all exomes sequenced (n=72), an average of 136 million high-quality 100bp paired-end reads were  
208 obtained. The total number of reads, mapping rate and coverage statistics for all sequencing runs are  
209 described in Supplementary Table 1. Mapping rates were between 94.2% and 99.9% with a median of  
210 98.9%. The median exome coverage rates were 97.2% (>30x) and 70.0% (>100x) of the exome.  
211 Sample HP\_4\_1 produced the lowest number of reads and subsequently had the lowest sequencing  
212 depth with 91.5% of the exome covered by >30x. Age group, sex, and counts for total SNVs, SNVs  
213 within a coding region (referred to CDS), and percentages of variants within a coding sequence region  
214 that correspond to known motifs for AID, ADAR, APOBEC3G and APOBEC3B are presented in  
215 Tables 1 and 2, and Supplementary Figure 3.

216

217 **Table 1:** Summary of total SNV counts, CDS SNV counts, and percentages of CDS variants that  
 218 correspond to motifs for AID, ADAR, APOBEC3G and APOBEC3B for Group 1 subjects, comprising  
 219 one blood and four saliva replicate datasets. Sex and age group is given for each healthy subject.

ID		HP_15	HP_16	HP_3	HP_10	HP_1	HP_4	HP_12	HP_20	HP_2	HP_5	HP_6	HP_22
Sex		M	M	F	F	M	M	F	F	M	M	F	F
Age Group		18-19	18-19	18-19	18-19	30-39	30-39	30-39	30-39	50-59	50-59	50-59	50-59
SNV	saliva-1	44227	44334	44913	44075	44411	44249	44531	44848	44121	44389	44256	44306
SNV	saliva-1C	44261	44408	45008	44196	44366	44354	44637	44914	44218	44444	44423	44419
SNV	saliva-2	44346	44509	44910	44222	44387	44291	44609	45051	44133	44304	44391	44356
SNV	saliva-2A	44305	44459	44934	44241	44436	44354	44616	44974	44170	44372	44414	44430
SNV	blood-3	44247	44338	44890	44124	44479	44297	44520	44963	44073	44340	44398	44324
CDS	saliva-1	20236	20309	20649	20084	20358	20326	20399	20605	20245	20359	20054	20233
CDS	saliva-1C	20280	20369	20651	20123	20358	20308	20418	20663	20306	20345	20131	20250
CDS	saliva-2	20293	20379	20641	20121	20380	20283	20408	20713	20275	20313	20136	20256
CDS	saliva-2A	20268	20398	20634	20126	20373	20298	20427	20655	20301	20327	20148	20277
CDS	blood-3	20255	20305	20620	20127	20415	20289	20389	20646	20255	20301	20136	20218
AID%	saliva-1	13.72	13.71	13.80	13.63	14.03	13.69	13.68	13.63	13.63	13.56	13.60	13.88
AID%	saliva-1C	13.67	13.71	13.82	13.64	14.00	13.67	13.67	13.69	13.66	13.62	13.59	13.84
AID%	saliva-2	13.73	13.72	13.83	13.66	14.03	13.69	13.69	13.66	13.68	13.57	13.61	13.88
AID%	saliva-2A	13.69	13.72	13.79	13.66	14.02	13.69	13.65	13.63	13.62	13.60	13.63	13.86
AID%	blood-3	13.72	13.69	13.81	13.60	14.02	13.70	13.68	13.62	13.68	13.55	13.60	13.90
ADAR%	saliva-1	15.40	15.43	15.67	15.85	15.57	15.54	15.45	15.64	15.56	15.66	15.45	15.64
ADAR%	saliva-1C	15.43	15.43	15.67	15.83	15.59	15.50	15.46	15.65	15.54	15.64	15.45	15.64
ADAR%	saliva-2	15.37	15.41	15.66	15.83	15.58	15.48	15.47	15.64	15.60	15.67	15.47	15.61
ADAR%	saliva-2A	15.42	15.39	15.69	15.81	15.59	15.50	15.46	15.62	15.57	15.66	15.47	15.58
ADAR%	blood-3	15.37	15.47	15.69	15.85	15.52	15.49	15.48	15.67	15.60	15.69	15.48	15.59
APOBEC3G%	saliva-1	16.95	16.71	16.90	16.68	16.96	16.92	17.13	16.90	16.76	16.90	16.74	16.65
APOBEC3G%	saliva-1C	16.94	16.75	16.89	16.72	16.91	16.92	17.16	16.85	16.77	16.90	16.71	16.65
APOBEC3G%	saliva-2	16.95	16.71	16.90	16.71	16.91	16.97	17.12	16.85	16.69	16.92	16.72	16.58
APOBEC3G%	saliva-2A	16.97	16.70	16.88	16.67	16.93	16.94	17.12	16.91	16.78	16.92	16.70	16.64
APOBEC3G%	blood-3	17.01	16.69	16.86	16.74	16.95	17.03	17.11	16.85	16.75	16.92	16.73	16.65
APOBEC3B%	saliva-1	4.06	4.13	3.99	4.09	4.13	4.02	4.17	4.06	3.99	4.13	4.26	4.14
APOBEC3B%	saliva-1C	4.05	4.12	4.00	4.11	4.13	4.03	4.12	4.03	3.98	4.10	4.26	4.13
APOBEC3B%	saliva-2	4.07	4.10	4.00	4.10	4.11	4.02	4.14	4.03	3.99	4.11	4.26	4.14
APOBEC3B%	saliva-2A	4.05	4.12	4.00	4.10	4.11	4.02	4.15	4.03	3.99	4.13	4.26	4.12
APOBEC3B%	blood-3	4.05	4.11	4.00	4.10	4.13	4.00	4.14	4.03	3.99	4.13	4.26	4.14

220

221

222

223 **Table 2:** Summary of total SNV counts, CDS SNV counts, and percentages of CDS variants that  
 224 correspond to the motifs for AID, ADAR, APOBEC3G and APOBEC3B for Group 2 subjects,  
 225 comprising saliva-1 samples. Sex and age group is given for each healthy subject.

ID		HP_14	HP_17	HP_11	HP_18	HP_9	HP_13	HP_7	HP_19	HP_8	HP_21	HP_23	HP_24
Sex		M	M	F	F	M	M	F	F	M	M	F	F
Age Group		18-19	18-19	18-19	18-19	30-39	30-39	30-39	30-39	50-59	50-59	50-59	50-59
SNV	saliva-1	44427	44632	45094	44540	44122	44643	44594	44491	44404	44425	44974	44258
CDS	saliva-1	20379	20480	20712	20431	20116	20380	20514	20390	20216	20452	20562	20290
AID%	saliva-1	13.67	13.70	13.45	13.61	13.70	13.67	13.84	13.55	13.76	13.86	13.89	13.64

	1												
ADAR%	saliva-1	15.52	15.42	15.52	15.55	15.66	15.64	15.81	15.82	15.50	15.49	15.49	15.44
APOBEC3G%	saliva-1	17.10	17.01	16.90	16.63	16.63	17.18	16.63	16.77	16.75	16.85	16.90	17.13
APOBEC3B%	saliva-1	3.93	4.11	4.04	4.11	4.21	3.92	3.95	3.97	3.97	4.09	4.11	3.97

226

## 227 SNV concordance between and within sample types

228 For Group 1 subjects (n=12), SNVs called in each sample were analyzed following the workflow  
 229 described in Figure 1. Variants shared between sample types were quantified (i.e. concordance between  
 230 saliva ‘1’, ‘1C’, ‘2’, ‘2A’, and blood ‘3’), with all sample types showing very high concordance  
 231 overall. Venn diagrams illustrate overlap of variant calls in all exome regions, as well as those located  
 232 within gene coding regions (CDS) between sample types and replicates for a representative volunteer  
 233 (HP\_1: Figure 2A and 2C). Overall, 96.1% of total variants in this volunteer were common to all  
 234 sample replicates and are referred to as *concordant* SNVs. Venn diagrams for all volunteers are  
 235 provided in Supplementary Figures 1 (all SNVs) and 2 (SNVs restricted to gene CDS). SNVs common  
 236 to 1, 2, 3, or 4 but not all 5 of the samples are referred to as *discordant* SNVs and were further  
 237 investigated. The percentage of concordant SNVs was slightly higher on average in the CDS (96.6%),  
 238 compared to those in all WES regions (95.8%) (Supplementary Table 2). As a measure of pairwise  
 239 similarity between samples, the number of *discordant* SNVs in common between sample pairs are  
 240 shown in Figure 2B, WES SNVs, and Figure 2D, CDS SNVs. Pairwise comparisons are categorized as:  
 241 biological and technical *blood-saliva replicates* (blood-3 & saliva-1, blood-3 & saliva-2, blood-3 &  
 242 saliva-1C, blood-3 & saliva-2A), biological and technical *saliva replicates* (saliva-1 & saliva-2, saliva-  
 243 1 & saliva-2A, saliva-1C & saliva-2, saliva-1C & saliva-2A), and *technical saliva replicates* (saliva-1  
 244 & saliva-1C, and saliva-2 & saliva-2A). The sample with lowest coverage (HP\_4 saliva-1) was  
 245 associated with lower pairwise overlap of discordant reads, however this difference was ameliorated  
 246 when analysis was restricted to only the coding region of genes. A statistical analysis of the pairwise  
 247 number of discordant SNVs in common (WES and CDS SNVs) showed no significant difference

248 between the pairwise comparisons of blood-saliva, saliva-saliva and technical saliva replicates  
249 (ANOVA;  $p > 0.05$ ).

250

251 Sequencing depth for concordant SNVs and discordant SNVs, averaged across 12 samples, is presented  
252 in Figure 3. Sequencing coverage distribution typically centered around 100x. Discordant SNVs have a  
253 higher density of low WES coverage ( $< 30x$ ). Depth analysis of individual samples are graphed in  
254 Supplementary Figure 4. Analysis of HP\_1 replicates revealed discordant SNVs failed one or more  
255 quality filters due to high strand bias  $> 10$ , (13% of 1762 discordant SNVs), low genotype quality  
256 (62%), high ratio of quality-filtered bases (8%), low depth (7%), or were not called as variants in one  
257 or more samples (54%). Overlap of concordant and discordant SNVs with dbSNP and gnomAD  
258 databases showed clear differences (Figure 3B). Concordant SNVs have a much higher overlap in  
259 dbSNP than discordant SNVs for both ‘all variants’ and ‘common’ variants. Using large-scale analysis  
260 of over 120 thousand exomes, the gnomAD database flags variants that do not pass certain quality  
261 filters. Of all concordant SNVs, 3% failed the gnomAD filters, however 31% of the discordant SNVs  
262 failed.

263

#### 264 **Candidate somatic SNV analysis**

265 Candidate somatic variants were identified using Strelka2 ‘tumor-normal’ methods, with blood and  
266 saliva sample replicates alternatingly used as ‘tumor’ and ‘normal’. There was no overlap between  
267 discordant SNVs identified in the germline and candidate somatic SNVs. Although all candidate  
268 somatic variants passed default filters, the quality of candidate somatic SNVs measured using the  
269 Strelka2 Empirical Variant Score (EVS) were all relatively low ( $EVS < 20$ ). EVS is a phred-scaled  
270 probability of the call being a false positive observation and is calculated from pre-trained random  
271 forest models and not hard cutoffs[31]. Low EVS scores are typically due to low minor allele  
272 frequencies, low mapping quality and low sequence coverage regions[31].

273 The mean number of somatic SNV candidates found in saliva was 149 (saliva=tumor, blood=normal),  
274 those found in blood was 121 (blood=tumor, saliva=normal). There was no correlation detected  
275 between the number of candidate somatic variants and the age of subjects. The average number of  
276 candidate saliva-blood somatic SNVs was 158, 141 and 148, and blood-saliva averages were 119, 131  
277 and 111 across the age groups 18-19, 30-39 and 50-59 respectively. A measure of technical noise is  
278 given by the number of candidate somatic variants found in biological and technical saliva replicates,  
279 which were on average 126 and 123 SNVs respectively (Supplementary Figure 5).

280 Of the candidate somatic SNVs identified, approximately 80% had a variant minor allele frequency  
281 <0.05. Applying this conventional filter reduced the average number of candidates per category (saliva,  
282 blood, technical replicates, biological replicates) to 34, 31, 23 and 25 respectively. A minimum depth  
283 filter of >30 for both ‘tumor’ and ‘normal’ samples further reduced average number of somatic  
284 candidates per category to 22, 19, 14 and 15 respectively.

285 After filtering of the candidate saliva SNVs that were not detected in all saliva replicates, 22 candidates  
286 remained with 21 of these found in more than one subject. Manual inspection using IGV suggests a  
287 false positive caused by incorrect mapping of soft clipped reads. With only a single blood sample per  
288 subject, equivalent filtering of candidate somatic SNVs found only in blood was not possible. The  
289 number of candidate somatic SNVs in blood was no larger than the level of technical and biological  
290 noise.

291

### 292 **Deaminase associated SNVs**

293 SNVs corresponding to known deaminase motifs (AID: WRC / GYW, ADAR: WA / TW,  
294 APOBEC3G: CC / GG, APOBEC3B: TCW / WGA) were identified within the coding region of genes  
295 (Tables 1 & 2, Supplementary Figure 3). Deaminase-associated SNVs at the genotype level were  
296 highly concordant and similar to the percentage concordance of all CDS SNVs: AID (96.1%), ADAR  
297 (97.0%), APOBEC3G (96.3%), APOBEC3B (96.2%) and CDS (96.6%) (Supplementary Table 2).

298 Strand bias and transition/transversion ratios were calculated for each deaminase and are summarized  
 299 below (Table 3). In total, approximately 50% of the CDS SNVs correspond to AID, ADAR,  
 300 APOBEC3G and APOBEC3B deaminase motifs. The intraclass correlation coefficient (ICC) was used  
 301 to quantify the reproducibility of the deaminase metric calculations. The ICC values range between  
 302 0.958 and 0.989 for all SNVs, deaminase motifs and associated metrics, revealing a very high  
 303 consistency/reproducibility among the five samples. ICC calculations for a range of additional metrics  
 304 are reported in Supplementary Table 3.

305

306

307 **Table 3:** Intraclass correlation coefficients (ICC) illustrating consistency between replicates in relation  
 308 to SNV calls and the existence of deaminase motifs, transition/transversion metrics (Ti/Tv), and strand  
 309 bias metrics. The variation *within* Group 1 replicates (n=12, 60 WES datasets) represents a measure of  
 310 technical reproducibility across blood and saliva samples. Variation *between* samples (n=24 Group 1  
 311 and 2 saliva-1 samples) represents a measure of biological variability.

	Mean <sup>a</sup>	Variance <i>within</i> replicates <sup>b</sup>	Variance <i>between</i> individuals <sup>c</sup>	ICC
Exome SNVs	44469	57.371	273.809	0.958
CDS SNVs	20366	25.0273	170.131	0.979
AID %	13.704	0.020	0.128	0.976
ADAR %	15.572	0.021	0.129	0.974
APOBEC3G %	16.862	0.026	0.169	0.978
APOBEC3B %	4.065	0.009	0.091	0.989
AID Ti/Tv %	74.166	0.081	0.550	0.979
AID C:G %	50.708	0.080	0.360	0.953
ADAR Ti/Tv %	77.908	0.071	0.386	0.968
ADAR A:T %	56.326	0.082	0.455	0.968
APOBEC3G Ti/Tv %	72.339	0.063	0.372	0.972
APOBEC3G C:G %	52.297	0.083	0.643	0.983
APOBEC3B Ti/Tv %	59.233	0.160	1.112	0.980
APOBEC3B C:G %	49.621	0.172	1.104	0.976

312 <sup>a</sup> Mean values from all saliva-1 samples (n=24, Group 1 and Group 2 subjects, 24 WES datasets)

313 <sup>b</sup> Average standard deviation from five replicates per volunteer (n=12, Group 1, 60 WES datasets)

314 <sup>c</sup> Standard deviation from all saliva-1 samples (n=24, Group 1 and Group 2 subjects, 24 WES datasets)

315

## 316 **Analysis of unmapped reads**

317 The average number of unmapped reads was larger in saliva (60 WES datasets, mean=2,372,300,  
318 98.4% mapping rate) than in blood (12 WES datasets, mean=334,182, 99.7% mapping rate),  
319 corresponding to a six fold higher unmapped rate in saliva (1.63% unmapped) compared to blood  
320 (0.27% unmapped). Overall, there is a 98.6% average mapping across all 72 samples and replicates  
321 (Supplementary Table 1). Quality statistics for unmapped reads are summarized at  
322 <https://jpmam1.github.io/MultiQC/>. Unmapped reads for volunteers HP\_1 and HP\_2 were extracted  
323 and aligned to the nr protein database. with read alignment rate to the NCBI nr database larger in saliva  
324 (41%) than in blood (33%). Reads that failed to align to NCBI nr were typically low quality.  
325 Unmapped reads derived from saliva, but not blood, were predominantly found to contain reads  
326 aligning to metagenomic species (Supplementary Figure 6).

327

## 328 **Discussion**

329 AID, APOBEC3G, APOBEC3B and ADAR deaminases are implicated in 30%-40% of clinically cu-  
330 rated SNPs in the OMIM database [13]. However, there is a paucity of research on deaminase-  
331 associated motifs in healthy subjects. Here, we have investigated deaminase-associated signatures in  
332 blood and saliva of healthy Caucasian subjects using whole exome sequencing. Our experimental de-  
333 sign provided a framework to quantify variance in all SNVs, SNVs within the coding sequence of  
334 genes, deaminase-associated coding SNVs, and provided measures of deaminase strand-bias and transi-  
335 tion/transversion in a highly robust and reproducible manner. Using different biological and technical  
336 sample replicates we explored differences between concordant and discordant SNV calls across the 24-  
337 subject cohort, showing strong intraclass correlation between sample replicates. No significant differ-  
338 ences in discordant SNV calls were detected in pairwise comparisons between sample types. The  
339 sources of discordant SNVs were investigated and were found to be associated with low read depth,



340 high strand bias, and low genotype quality. Analysis of putative somatic variants showed no conclusive  
341 evidence of somatic mutation when comparing blood and saliva samples. On average, approximately  
342 2% of reads failed to align to the human genome, with reads derived from saliva samples primarily re-  
343 lated to metagenomic taxa associated with the oral microbiome [41,42]. Here, we establish that saliva  
344 and blood are both appropriate sources of DNA for WES analyses, with no detected difference in abil-  
345 ity to resolve SNVs and deaminase-associated signatures and metrics.

346

347 A key component of the experimental design in this study (Figure 1) was the capacity for comparisons  
348 between biological and technical replicates. The replicate extraction of Saliva-1 DNA, and replicate  
349 library preparation of Saliva-2, provides a measure of lab-based technical variation. Our study showed  
350 that the differences between blood and saliva, and between biological saliva replicates were very small  
351 and no larger than the level of noise of the technical replicates. Discordant SNV calls are  
352 predominantly in low coverage and/or low-quality regions of the exome. These discordant SNVs were  
353 less likely to be found in dbSNP and were dramatically enriched for known problematic SNV calls in  
354 gnomAD. These results indicate the filtering for SNVs that fail gnomAD quality analysis would  
355 improve overall reproducibility of WES SNV analysis. By filtering these failed gnomAD variants, only  
356 3% of concordant SNVs are eliminated, but 30% of discordant SNVs are removed. These results may  
357 advise filtering strategies in future studies.

358

359 Previous research has shown high quality DNA can be obtained from saliva and blood, with results  
360 from whole exome sequencing found to be comparable for specific applications [25,26,43]. Due to oral  
361 microbiome and off-target capture, the overall unmapped rate of saliva in this study was six fold higher  
362 than that of blood (1.6% vs 0.3%), providing sufficient unmapped reads to perform a limited metage-  
363 nomic analysis. Given the importance of the relationship between the microbiome and immunity [44],



364 the oral microbiome information provided from off-target saliva exome capture may prove useful for a  
365 variety of applications in future studies (e.g. Kidd et al. [25]).

366

367 Accumulation of a small number of ‘age-related’ (pre- or non-cancerous) somatic mutations has been  
368 reported in several recent studies [16,45,46]. Despite our comprehensive analysis of WES data for  
369 evidence of somatic mutations across different ages and sexes, we were unable to unambiguously  
370 detect somatic mutations by comparing blood and saliva in these healthy individuals. Our use of  
371 biological and technical saliva replicates revealed similar numbers of candidate somatic SNVs in both  
372 technical and biological replicates for all subjects. This indicates a high level of noise and coincides  
373 with recent analyses of false-positive variant calls [47]. Bespoke somatic SNV detection approaches are  
374 evidently required to identify somatic SNVs in healthy subjects, using more advanced sequencing  
375 techniques, different cell types and more sophisticated bioinformatics [16].

376

377 There are many challenges in accurately resolving germline and somatic SNVs. Sequencing and bioin-  
378 formatics artefacts are known to result in incorrect SNV calls, with numerous studies investigating the  
379 effects of exome capture kits, sequencing platform, and bioinformatics software on the ability to accu-  
380 rately detect SNVs [48,49]. As reported by others, performance of pipelines according to a ‘gold stan-  
381 dard’ (such as “genome in a bottle”) does not necessarily indicate performance on ‘real world’ datasets  
382 [50]. In this study, the use of sample replicates enabled us to quantify noise and evaluate the reproduc-  
383 ibility of SNV calls, to identify discordant germline SNVs as potential false positives and eliminate  
384 false-positive somatic SNVs. Reducing false-positive SNVs is necessary to accurately resolve deami-  
385 nase-associated SNV profiles and for understanding the implications of deaminase signatures in health  
386 and disease.

387

388 Deaminase mutagenesis is associated with an increasing number of viral infections and cancer types  
389 [6,19,51–58]. With development of more advanced sequencing technologies, we will be able to detect  
390 evidence of deaminase activity with greater accuracy and examine changes over time [1,59,60]. In ad-  
391 dition to the well-characterized effects of deaminases on genome stability in cancer, and more recently  
392 in precancerous conditions [55,61,62], deaminases have emerged as a driving factor in many human  
393 SNPs [13]. Despite the limitations of 24 individuals, and all having Caucasian ancestry, this study en-  
394 abled us to investigate candidate somatic SNVs and provided us with a robust and reliable measure of  
395 deaminase-associated germline variants in healthy subjects.

396

## 397 **Conclusions**

398 A large proportion of disease-associated germline variants are linked to deaminase activity. We have  
399 established that saliva and blood are appropriate sources of DNA for whole exome sequencing, with no  
400 difference in ability to resolve deaminase-associated metrics. Deaminase-associated mutations are  
401 important in pre-cancerous conditions, and in cancer, however no somatic SNVs were identified when  
402 comparing blood and saliva of healthy individuals. Investigation into the implications of deaminase  
403 activity on genome stability in healthy individuals will required more technically advanced approaches.

404

## 405 **Abbreviations**

406 **ADAR:** Adenosine Deaminase Acting on RNA;

407 **AID:** activation induced cytidine deaminase, a APOBEC family member;

408 **APOBEC family:** generic abbreviation for the deoxyribonucleic acid deaminase family (APOBECs 1,2,4 and  
409 3A/B/C/D/F/G/H);

410 **CDS:** Coding sequence.

411 **Deaminase:** zinc-containing catalytic domain in ADAR and AID/APOBEC enzymes;

412 **HP:** Healthy Population (or Person);

413 **ICC:** Intraclass Correlation Coefficient;  
414 **R:** Adenosine (A) or Guanine (G), purines;  
415 **SD:** standard deviation;  
416 **SNP: single nucleotide polymorphism;**  
417 **SNV:** single nucleotide variant;  
418 **TSM:** targeted somatic mutations;  
419 **W:** weak base pair involving A or T;  
420 **Y:** pyrimidines T or C.;  
421 **WES:** whole exome sequencing.

422

423

## 424 **Declarations**

425

### 426 **Ethics approval and consent to participate**

427 This project was approved by the Monash Health Human Research Ethics Committee (16281L: “A  
428 study to measure the Targeted Somatic Mutation (TSM) test platform performance characteristics and  
429 evaluate its suitability for clinical use”). All study subjects signed written informed consent forms  
430 which were approved by the Ethics Committee.

431

### 432 **Consent for publication:**

433 “Not applicable”

434

### 435 **Availability of data and material**

436 The data are not publicly available due to information that could compromise research participant  
437 privacy and consent. The data that support the findings of this study are available on reasonable request  
438 from the corresponding author NEH.

439

440 **Competing interests**

441 All authors declare that they have no competing interests.

442

443 **Funding**

444 The work was fully funded by GMDx Group Ltd (Melbourne, Australia), as a part of a GMDx Group  
445 translational research program.

446

447 **Authors' contributions**

448 RAL conceived the study, NEH and RAL designed the study with CF and RB. NEH and JM analyzed  
449 and interpreted the genomic data and wrote the manuscript. PR and RB were involved in implementing  
450 the clinical trial. RAL and EJS contributed to the writing of the manuscript. All authors read and  
451 approved the final manuscript.

452

453 **Acknowledgements**

454 We thank Christopher Pendlebury and Richard Rendell from Applied Precision Medicine for TSM  
455 computational platform development and implementation. The authors also thank Trevor Wilson, Niro  
456 Pathirage, Roxane Legaie and Wishva Herath from the Monash Health Translation Precinct (MHTP)  
457 Medical Genomics Facility for exome sequencing and data processing and Margaret Smith from  
458 smartDNA for DNA extraction. Lastly, we wish to acknowledge the contributions of the study  
459 volunteers who provided samples.

460

461

462 **References:**

- 463 1. Casellas R, Basu U, Yewdell WT, Chaudhuri J, Robbiani DF, Noia JMD. Mutations, kataegis and  
464 translocations in B cells: understanding AID promiscuous activity. *Nat Rev Immunol.* 2016;16:164–76.
- 465 2. Conticello SG. The AID/APOBEC family of nucleic acid mutators. *Genome Biol.* 2008;9:229.
- 466 3. Harris RS. Cancer mutation signatures, DNA damage mechanisms, and potential clinical implica-  
467 tions. *Genome Med.* 2013;5:87.
- 468 4. Smith HC, Bennett RP, Kizilyer A, McDougall WM, Prohaska KM. Functions and regulation of the  
469 APOBEC family of proteins. *Semin Cell Dev Biol.* 2012;23:258–68.
- 470 5. Tubbs A, Nussenzweig A. Endogenous DNA Damage as a Source of Genomic Instability in Cancer.  
471 *Cell.* 2017;168:644–56.
- 472 6. Burns MB, Temiz NA, Harris RS. Evidence for APOBEC3B mutagenesis in multiple human can-  
473 cers. *Nat Genet.* 2013;45:977–83.
- 474 7. Jarvis MC, Ebrahimi D, Temiz NA, Harris RS. Mutation Signatures Including APOBEC in Cancer  
475 Cell Lines. *JNCI Cancer Spectr.* 2018; 2. Available from:  
476 <https://academic.oup.com/jncics/article/2/1/pky002/4942295>
- 477 8. Nikkilä J, Kumar R, Campbell J, Brandsma I, Pemberton HN, Wallberg F, et al. Elevated  
478 APOBEC3B expression drives a kataegic-like mutation signature and replication stress-related thera-  
479 peutic vulnerabilities in p53-defective cells. *Br J Cancer.* 2017;117:113–23.
- 480 9. Swanton C, McGranahan N, Starrett GJ, Harris RS. APOBEC Enzymes: Mutagenic Fuel for Cancer  
481 Evolution and Heterogeneity. *Cancer Discov.* 2015;5:704–12.
- 482 10. Steele EJ, Lindley RA. ADAR and APOBEC editing signatures in viral RNA during acute-phase  
483 Innate Immune responses of the host-parasite relationship to Flaviviruses. *Res Rep.* 2018; Available  
484 from: <http://www.companyofscientists.com/index.php/rr/article/view/80>
- 485 11. Lodato MA, Rodin RE, Bohrson CL, Coulter ME, Barton AR, Kwon M, et al. Aging and  
486 neurodegeneration are associated with increased mutations in single human neurons. *Science.*  
487 2017;eaao4426.
- 488 12. Spira A, Yurgelun MB, Alexandrov L, Rao A, Bejar R, Polyak K, et al. Precancer Atlas to Drive  
489 Precision Prevention Trials. *Cancer Res.* 2017;77:1510–41.
- 490 13. Lindley RA, Hall NE. APOBEC and ADAR deaminases may cause many single nucleotide poly-  
491 morphisms curated in the OMIM database. *Mutat Res Mol Mech Mutagen.* 2018;810:33–8.
- 492 14. Matthews AJ, Husain S, Chaudhuri J. Binding of AID to DNA Does Not Correlate with Mutator  
493 Activity. *J Immunol.* 2014;1400433.
- 494 15. Venkatesan S, Rosenthal R, Kanu N, McGranahan N, Bartek J, Quezada SA, et al. Perspective:  
495 APOBEC mutagenesis in drug resistance and immune escape in HIV and cancer evolution. *Ann Oncol.*  
496 2018;29:563–72.
- 497 16. Dou Y, Gold HD, Luquette LJ, Park PJ. Detecting Somatic Mutations in Normal Cells. *Trends*  
498 *Genet.* 2018;34:545–57.

- 499 17. Krishnan A, Iyer LM, Holland SJ, Boehm T, Aravind L. Diversification of AID/APOBEC-like  
500 deaminases in metazoa: multiplicity of clades and widespread roles in immunity. *Proc Natl Acad Sci*.  
501 2018;201720897.
- 502 18. Siriwardena SU, Chen K, Bhagwat AS. Functions and Malfunctions of Mammalian DNA-Cytosine  
503 Deaminases. *Chem Rev*. 2016;116:12688–710.
- 504 19. Lindley RA. The importance of codon context for understanding the Ig-like somatic hypermutation  
505 strand-biased patterns in TP53 mutations in breast cancer. *Cancer Genet*. 2013;206:222–6.
- 506 20. Refsland EW, Harris RS. The APOBEC3 Family of Retroelement Restriction Factors. *Curr Top*  
507 *Microbiol Immunol*. 2013;371:1–27.
- 508 21. Hansen T v O, Simonsen MK, Nielsen FC, Hundrup YA. Collection of Blood, Saliva, and Buccal  
509 Cell Samples in a Pilot Study on the Danish Nurse Cohort: Comparison of the Response Rate and  
510 Quality of Genomic DNA. *Cancer Epidemiol Prev Biomark*. 2007;16:2072–6.
- 511 22. Jobling MA, Gill P. Encoded evidence: DNA in forensic analysis. *Nat Rev Genet*. 2004;5:739–51.
- 512 23. Quinque D, Kittler R, Kayser M, Stoneking M, Nasidze I. Evaluation of saliva as a source of hu-  
513 man DNA for population and association studies. *Anal Biochem*. 2006;353:272–7.
- 514 24. Zaaijer S, Gordon A, Speyer D, Piccone R, Groen SC, Erlich Y. Rapid re-identification of human  
515 samples using portable DNA sequencing. *eLife*. 2017. Available from:  
516 <https://elifesciences.org/articles/27798>
- 517 25. Kidd JM, Sharpton TJ, Bobo D, Norman PJ, Martin AR, Carpenter ML, et al. Exome capture from  
518 saliva produces high quality genomic and metagenomic data. *BMC Genomics*. 2014;15:262.
- 519 26. Patel ZH, Kottyan LC, Lazaro S, Williams MS, Ledbetter DH, Tromp G, et al. The struggle to find  
520 reliable results in exome sequencing data: filtering out Mendelian errors. *Front Genet*. 2014; 5. Availa-  
521 ble from: <https://www.frontiersin.org/articles/10.3389/fgene.2014.00016/full>
- 522 27. Bruinsma FJ, Joo JE, Wong EM, Giles GG, Southey MC. The utility of DNA extracted from saliva  
523 for genome-wide molecular research platforms. *BMC Res Notes*. 2018;11:8.
- 524 28. Andrews. FastQC A Quality Control tool for High Throughput Sequence Data. 2010. Available  
525 from: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- 526 29. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
527 *EMBnet.journal*. 2011;17:10–2.
- 528 30. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
529 *ArXiv13033997 Q-Bio*. 2013; Available from: <http://arxiv.org/abs/1303.3997>
- 530 31. Kim S, Scheffler K, Halpern AL, Bekritsky MA, Noh E, Källberg M, et al. Strelka2: fast and accu-  
531 rate calling of germline and somatic variants. *Nat Methods*. 2018;15:591–4.
- 532 32. Hinrichs AS, Karolchik D, Baertsch R, Barber GP, Bejerano G, Clawson H, et al. The UCSC Ge-  
533 nome Browser Database: update 2006. *Nucleic Acids Res*. 2006;34:D590–8.

- 534 33. Zerbino DR, Achuthan P, Akanni W, Amode MR, Barrell D, Bhai J, et al. Ensembl 2018. *Nucleic*  
535 *Acids Res.* 2018;46:D754–61.
- 536 34. Sherry ST, Ward M-H, Kholodov M, Baker J, Phan L, Smigielski EM, et al. dbSNP: the NCBI da-  
537 tabase of genetic variation. *Nucleic Acids Res.* 2001;29:308–11.
- 538 35. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-  
539 coding genetic variation in 60,706 humans. *Nature.* 2016;536:285–91.
- 540 36. Ewels P, Magnusson M, Lundin S, Käller M. MultiQC: summarize analysis results for multiple  
541 tools and samples in a single report. *Bioinformatics.* 2016;32:3047–8.
- 542 37. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat Meth-*  
543 *ods.* 2015;12:59–60.
- 544 38. Beier S, Tappu R, Huson DH. Functional Analysis in Metagenomics Using MEGAN 6. In: Charles  
545 TC, Liles MR, Sessitsch A, editors. *Funct Metagenomics Tools Appl.* Cham: Springer International  
546 Publishing; 2017. p. 65–74. Available from: [https://doi.org/10.1007/978-3-319-61510-3\\_4](https://doi.org/10.1007/978-3-319-61510-3_4)
- 547 39. Huson DH, Beier S, Flade I, Górska A, El-Hadidi M, Mitra S, et al. MEGAN Community Edition -  
548 Interactive Exploration and Analysis of Large-Scale Microbiome Sequencing Data. *PLOS Comput*  
549 *Biol.* 2016;12:e1004957.
- 550 40. Shrout PE, Fleiss JL. Intraclass correlations: Uses in assessing rater reliability. *Psychol Bull.*  
551 1979;86:420–8.
- 552 41. Macovei L, McCafferty J, Chen T, Teles F, Hasturk H, Paster BJ, et al. The hidden  
553 ‘mycobacteriome’ of the human healthy oral cavity and upper respiratory tract. *J Oral Microbiol.*  
554 2015;7. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4329316/>
- 555 42. Verma D, Garg PK, Dubey AK. Insights into the human oral microbiome. *Arch Microbiol.*  
556 2018;200:525–40.
- 557 43. Zhu Q, Hu Q, Shepherd L, Wang J, Wei L, Morrison CD, et al. The Impact of DNA Input Amount  
558 and DNA source on the Performance of Whole-Exome Sequencing in Cancer Epidemiology. *Cancer*  
559 *Epidemiol Prev Biomark.* 2015;cebp.0205.2015.
- 560 44. Ost KS, Round JL. Communication Between the Microbiota and Mammalian Immunity. *Annu Rev*  
561 *Microbiol.* 2018; Available from: [http://www.annualreviews.org/doi/10.1146/annurev-micro-090817-](http://www.annualreviews.org/doi/10.1146/annurev-micro-090817-062307)  
562 [062307](http://www.annualreviews.org/doi/10.1146/annurev-micro-090817-062307)
- 563 45. Forsberg LA, Gisselsson D, Dumanski JP. Mosaicism in health and disease — clones picking up  
564 speed. *Nat Rev Genet.* 2017;18:128–42.
- 565 46. Vattathil S, Scheet P. Extensive Hidden Genomic Mosaicism Revealed in Normal Tissue. *Am J*  
566 *Hum Genet.* 2016;98:571–8.
- 567 47. Potapov V, Ong JL. Examining Sources of Error in PCR by Single-Molecule Sequencing. *PLOS*  
568 *ONE.* 2017;12:e0169774.
- 569 48. Hwang S, Kim E, Lee I, Marcotte EM. Systematic comparison of variant calling pipelines using  
570 gold standard personal exome variants. *Sci Rep.* 2015;5:17875.



- 571 49. Shigemizu D, Momozawa Y, Abe T, Morizono T, Boroevich KA, Takata S, et al. Performance  
572 comparison of four commercial human whole-exome capture platforms. *Sci Rep.* 2015;5:12742.
- 573 50. Bohnert R, Vivas S, Jansen G. Comprehensive benchmarking of SNV callers for highly admixed  
574 tumor data. *PLOS ONE.* 2017;12:e0186175.
- 575 51. Beale RCL, Petersen-Mahrt SK, Watt IN, Harris RS, Rada C, Neuberger MS. Comparison of the  
576 Differential Context-dependence of DNA Deamination by APOBEC Enzymes: Correlation with Muta-  
577 tion Spectra in Vivo. *J Mol Biol.* 2004;337:585–96.
- 578 52. Burns MB, Lackey L, Carpenter MA, Rathore A, Land AM, Leonard B, et al. APOBEC3B is an  
579 enzymatic source of mutation in breast cancer. *Nature.* 2013;494:366–70.
- 580 53. Gallo A, Locatelli F. ADARs: allies or enemies? The importance of A-to-I RNA editing in human  
581 disease: from cancer to HIV-1. *Biol Rev.* 2012;87:95–110.
- 582 54. Dominissini D, Moshitch-Moshkovitz S, Amariglio N, Rechavi G. Adenosine-to-inosine RNA edit-  
583 ing meets cancer. *Carcinogenesis.* 2011;32:1569–77.
- 584 55. Lindley RA, Humbert P, Larner C, Akmeemana EH, Pendlebury CRR. Association between target-  
585 ed somatic mutation (TSM) signatures and HGS-OvCa progression. *Cancer Med.* 2016;5:2629–40.
- 586 56. Roberts SA, Lawrence MS, Klimczak LJ, Grimm SA, Fargo D, Stojanov P, et al. An APOBEC  
587 cytidine deaminase mutagenesis pattern is widespread in human cancers. *Nat Genet.* 2013;45:970–6.
- 588 57. Roberts SA, Gordenin DA. Hypermutation in human cancer genomes: footprints and mechanisms.  
589 *Nat Rev Cancer.* 2014;14:786–800.
- 590 58. Taylor BJ, Nik-Zainal S, Wu YL, Stebbings LA, Raine K, Campbell PJ, et al. DNA deaminases  
591 induce break-associated mutation showers with implication of APOBEC3B and 3A in breast cancer  
592 kataegis. *eLife.* 2013;2:e00534.
- 593 59. Martincorena I, Roshan A, Gerstung M, Ellis P, Van Loo P, McLaren S, et al. Tumor evolution.  
594 High burden and pervasive positive selection of somatic mutations in normal human skin. *Science.*  
595 2015;348:880–6.
- 596 60. Rogozin IB, Pavlov YI, Goncarenco A, De S, Lada AG, Poliakov E, et al. Mutational signatures  
597 and mutable motifs in cancer genomes. *Brief Bioinform.* 2017;
- 598 61. Chan THM, Lin CH, Qi L, Fei J, Li Y, Yong KJ, et al. A disrupted RNA editing balance mediated  
599 by ADARs (Adenosine DeAminases that act on RNA) in human hepatocellular carcinoma. *Gut.*  
600 2014;63:832–43.
- 601 62. Leonard B, Hart SN, Burns MB, Carpenter MA, Temiz NA, Rathore A, et al. APOBEC3B  
602 Upregulation and Genomic Mutation Patterns in Serous Ovarian Carcinoma. *Cancer Res.*  
603 2013;73:7222–31.
- 604

## 605 **Figure Legends**

### 606 **Figure 1. Experimental design**



607 (A) Flow diagram of whole exome sequencing (WES) pipeline. DNA was extracted from two  
608 saliva samples “1” and “2” and a blood sample “3” to make libraries for whole exome sequencing.  
609 Lib-1C is made from a separate DNA extraction of Sample 1, and Lib-2A is a separate library  
610 made from the same Saliva-2 DNA extraction. Blood samples have only one DNA extraction and  
611 one library preparation. (B) WES data processing pipeline. Reads were aligned to the reference  
612 human genome, SNVs were called in mapped reads and variants associated with deaminase  
613 motifs were quantified. Unmapped reads were aligned to the NCBI ‘non-redundant’ (nr) to  
614 establish the taxonomic sources of reads.

615

616 **Figure 2: Relationship of SNV calls between and among sample replicates.**

617 (A) Venn diagram of all called variants by ID, (B) pairwise sample comparisons of all discordant SNVs  
618 in common for each WES dataset pair, (C) Venn diagram of CDS variants, (D) pairwise sample  
619 comparisons of all discordant CDS SNVs in common for each WES dataset pair. Blood-saliva pairwise  
620 comparisons are in shades of red. Saliva-1-saliva-2 comparisons are in shades of blue, and technical  
621 saliva replicates are in green.

622

623 **Figure 3: Sequencing depth and database overlap of concordant and discordant SNVs.**

624 (A) Combined depth profiles of concordant SNVs across five samples types compared to discordant  
625 SNVs, and (B) the rate of concordant and discordant SNVs belonging to each variant database (n=12,  
626 mean  $\pm$  95% CI).

627

628

629 **Supplementary Figure Legends**

630 **Supplementary Figure 1: Relationships of SNV calls between and among sample replicates.**

631 Venn diagrams of exome SNV calls for five replicates for twelve healthy subjects. Saliva-1, and saliva-  
632 2 are replicate saliva samples. Saliva-1C and saliva-2A are technical replicates of saliva-1 and saliva-2.  
633 SNVs called in all five samples are termed *concordant* SNVs, and those not in all samples are termed  
634 *discordant* SNVs. Venn diagrams were generated online at  
635 <http://bioinformatics.psb.ugent.be/webtools/Venn/>.

636

### 637 **Supplementary Figure 2: Relationships of CDS SNV calls between and among sample replicates.**

638 Venn diagrams of coding sequence (CDS) SNV calls for five replicates for twelve healthy subjects.

639

640

### 641 **Supplementary Figure 3: Calculated metrics for all sample replicates**

642 Bar graphs of 72 exome data sets across 24 healthy individuals. Bars are colored according to the  
643 biological or technical replicate. Metrics presented are:

644 **SNV**, total number SNV calls,

645 **CDS**, number of SNVs in the coding regions,

646 **AID%**, percentage of CDS SNVs matching the AID deaminase motif  $\text{WRC} / \text{GYW}$ ,

647 **ADAR%**, percentage of CDS SNVs matching the ADAR deaminase motif  $\text{WA} / \text{TW}$ ,

648 **APOBEC3G%** percentage of CDS SNVs matching the APOBEC3G deaminase motif  $\text{CC} / \text{GG}$ , and

649 **APOBEC3B%**, percentage of CDS SNVs matching the APOBEC3B deaminase motif  $\text{TCW} / \text{WGA}$

650 where  $W=A$  or  $T$ ,  $Y=C$  or  $T$ ,  $R=A$  or  $G$ .

651

### 652 **Supplementary Figure 4: Depth profiles of concordant and discordant blood and saliva replicate** 653 **SNVs.**

654 Sequencing depth of five sample types for concordant SNVs compared to discordant SNVs (n=12).

655

656

657 **Supplementary Figure 5: Distribution and density of candidate somatic SNVs**

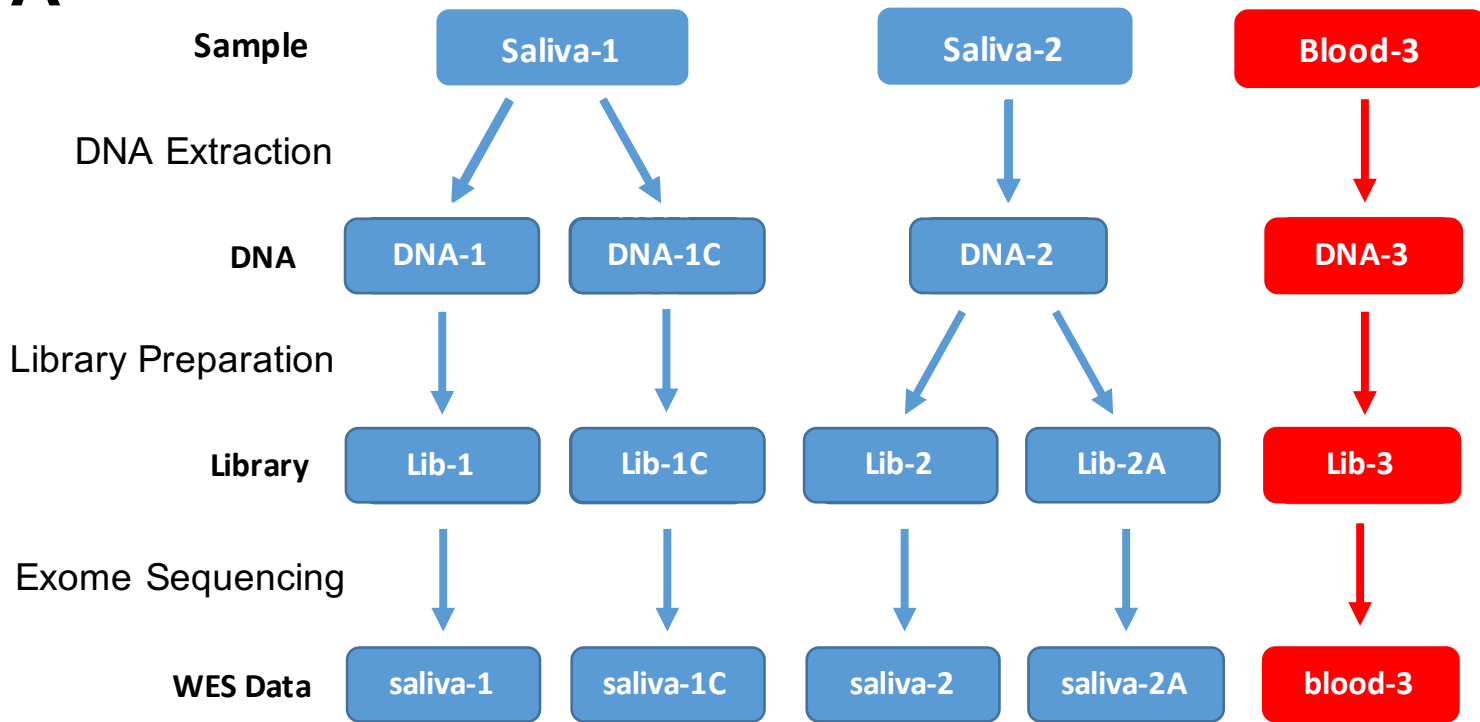
658 Number of candidate SNVs per sample, grouped by “tumor-normal” comparison type. For ‘blood vs  
659 saliva’, blood was treated as the *tumor* sample and saliva as *normal*; for “saliva vs blood”, saliva was  
660 treated as the *tumor* sample and blood as *normal*; for "saliva biological replicates", saliva-1 samples  
661 were compared against saliva-2 samples, and vice versa; and for "saliva technical replicates", saliva 1  
662 and saliva-1C samples were compared, and saliva-2 and saliva-2A samples were compared. Boxplots  
663 illustrate the median and interquartile range (IQR), with outliers defined as  $1.5 \times \text{IQR}$ . The distribution  
664 for each grouping is shown above each boxplot. The number of candidate SNVs detected between  
665 technical replicates demonstrates a high level of noise across all candidate somatic SNVs. Filtering and  
666 analysis suggests all candidate SNVs are likely false positives.

667

668 **Supplementary Figure 6: Metagenomic analysis of unmapped reads**

669 Alignment of unmapped reads for subjects HP\_1 and HP\_2 to the NCBI nr protein database  
670 representing the sources of off-target WES DNA. The majority of unmapped reads derived from saliva  
671 corresponded to prokaryotic organisms associated with the oral microbiome. Unmapped reads derived  
672 from blood samples were either low quality or predominantly mapped to viral DNA (Phi-X spike-in).

673

**A****B**