

1 High-throughput Interpretation of Killer-cell Immunoglobulin-like Receptor Short-read
2 Sequencing Data with PING

3

4 Wesley M. Marin¹, Ravi Dandekar¹, Danillo G. Augusto¹, Tasneem Yusufali¹, Bianca Heyn², Jan
5 Hofmann³, Vinzenz Lange², Jürgen Sauter³, Paul J. Norman⁴, Jill A. Hollenbach^{1*}

6

7 ¹UCSF Weill Institute for Neurosciences, Department of Neurology, University of California,
8 San Francisco, San Francisco, CA 94158, United States

9 ²DKMS Life Science Lab, Dresden 01069, Germany

10 ³DKMS, Tübingen 72072, Germany

11 ⁴Division of Biomedical Informatics and Personalized Medicine, and Department of
12 Immunology and Microbiology, University of Colorado Anschutz Medical, Aurora, CO 80047,
13 United States

14

15 * Corresponding author

16 Email: jill.hollenbach@ucsf.edu

17

18 **Abstract**

19 The *killer-cell immunoglobulin-like receptor* (*KIR*) complex on chromosome 19 encodes
20 receptors that modulate the activity of natural killer cells, and variation in these genes has been
21 linked to infectious and autoimmune disease, as well as having bearing on pregnancy and
22 transplant outcomes. The medical relevance and high variability of *KIR* genes makes short-read
23 sequencing an attractive technology for interrogating the region, providing a high-throughput,
24 high-fidelity sequencing method that is cost-effective. However, because this gene complex is
25 characterized by extensive nucleotide polymorphism, structural variation including gene fusions
26 and deletions, and a high level of homology between genes, its interrogation at high resolution
27 has been thwarted by bioinformatic challenges, with most studies limited to examining presence
28 or absence of specific genes. Here, we present the PING (Pushing Immunogenetics to the Next
29 Generation) pipeline, which incorporates empirical data, novel alignment strategies and a custom
30 alignment processing workflow to enable high-throughput *KIR* sequence analysis from short-
31 read data. PING provides *KIR* gene copy number classification functionality for all *KIR* genes
32 through use of a comprehensive alignment reference. The gene copy number determined per
33 individual enables an innovative genotype determination workflow using genotype-matched
34 references. Together, these methods address the challenges imposed by the structural complexity
35 and overall homology of the *KIR* complex. To determine copy number and genotype
36 determination accuracy, we applied PING to European and African validation cohorts and a
37 synthetic dataset. PING demonstrated exceptional copy number determination performance
38 across all datasets and robust genotype determination performance. Finally, an investigation into
39 discordant genotypes for the synthetic dataset provides insight into misaligned reads, advancing
40 our understanding in interpretation of short-read sequencing data in complex genomic regions.

41 PING promises to support a new era of studies of KIR polymorphism, delivering high-resolution
42 *KIR* genotypes that are highly accurate, enabling high-quality, high-throughput *KIR* genotyping
43 for disease and population studies.

44

45 **Author summary**

46 Killer cell immunoglobulin-like receptors (KIR) serve a critical role in regulating natural killer
47 cell function. They are encoded by highly polymorphic genes within a complex genomic region
48 that has proven difficult to interrogate owing to structural variation and extensive sequence
49 homology. While methods for sequencing *KIR* genes have matured, there is a lack of bioinformatic
50 support to accurately interpret *KIR* short-read sequencing data. The extensive structural variation
51 of *KIR*, both the small-scale nucleotide insertions and deletions and the large-scale gene
52 duplications and deletions, coupled with the extensive sequence similarity among *KIR* genes
53 presents considerable challenges to bioinformatic analyses. PING addressed these issues through
54 a highly-dynamic alignment workflow, which constructs individualized references that reflect the
55 determined copy number and genotype makeup of a sample. This alignment workflow is enabled
56 by a custom alignment processing pipeline, which scaffolds reads aligned to all reference
57 sequences from the same gene into an overall gene alignment, enabling processing of these
58 alignments as if a single reference sequence was used regardless of the number of sequences or of
59 any insertions or deletions present in the component sequences. Together, these methods provide
60 a novel and robust workflow for the accurate interpretation of *KIR* short-read sequencing data.

61

62 **Introduction**

63 The *killer cell immunoglobulin-like receptor (KIR)* complex, located in human chromosomal
64 region 19q13.42, encodes receptors expressed on the surface of natural killer (NK) cells (1) and a
65 subtype of T-cells (2). KIRs interact with their cognate HLA class I ligands to educate NK cells
66 and modulate their cytotoxicity (3–5). *KIR* genes exhibit presence and absence polymorphism and
67 gene content variation that has been implicated in numerous immune-mediated and infectious
68 diseases (6–11). In addition, careful consideration of *KIR* gene content haplotypes for allogeneic
69 transplantation has been shown to improve outcomes for acute myelogenous leukemia patients
70 (12–17). Whereas evidence for the relevance of *KIR* variation in health and disease is mounting,
71 analysis of the *KIR* family at allelic resolution has been thwarted by the complexity of the region.
72

73 The *KIR* complex evolved rapidly through recombination and gene duplication events, and in
74 humans this has resulted in a gene-content variable cluster of 13 genes and 2 pseudogenes (18–
75 20). Variation in *KIR* genes is characterized by extensive nucleotide polymorphisms, with 1110
76 alleles described to date (21). The *KIR* complex is also characterized by large-scale structural
77 variation, including gene fusions, duplications and deletions (22,23). *KIR* haplotypes exhibit gene
78 content variation at extraordinary levels, generating hundreds of observed haplotype structures
79 (20,24–26).

80
81 The high variability of *KIR* makes short-read sequencing an attractive technology for interrogating
82 the region, providing a high-throughput, high-fidelity and cost-effective sequencing method (27).
83 Whereas the *KIR* region is relatively small, between 70-270Kbp (28), the overall sequence
84 similarity among genes, structural variability of the region, and sequence polymorphism present

85 major obstacles to bioinformatics workflows. The high potential for read misalignments
86 significantly confounds interpretation of the region in modern large-scale sequencing studies.

87

88 Previously, we introduced a laboratory method for targeted sequencing of the *KIR* gene complex
89 (27), but the associated prototype bioinformatic pipeline for sequence interpretation presented
90 significant workload barriers for high-throughput studies. For example, the copy number
91 determination workflow was unable to differentiate *KIR2DL2* from *KIR2DL3*, which are sets of
92 highly similar allelic groups of the *KIR2DL23* gene (29), and the resolution of *KIR2DS1* and
93 *KIR2DL1* was less precise than desired due to read misalignments caused by the close similarity
94 of these two genes. Additionally, a high frequency of unresolved genotypes (not matching any
95 described allele sequence) necessitated subsequent interpretation by a user with domain expertise.
96 In spite of these challenges, the prototype pipeline has provided insight into *KIR* genotyping
97 methods development (30,31), the role of *KIR* sequence variants in immune dysfunction (17,32),
98 and *KIR* evolutionary analyses (33).

99

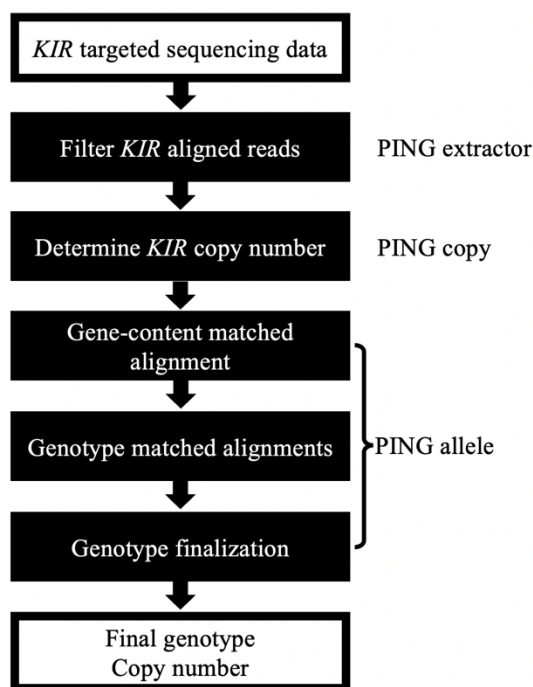
100 Here, we present a comprehensive *KIR* sequence interpretation workflow, termed PING (Pushing
101 Immunogenetics to the Next Generation), which builds on our early work by incorporating
102 empirical optimizations derived from sequencing thousands of samples, in addition to novel
103 alignment strategies to address issues with read misalignments. These innovations enable for the
104 first time highly-automated, high-throughput *KIR* sequence analysis from short-read sequencing
105 data, and, importantly, largely obviating the need for user expertise in the *KIR* system.

106

107 **Materials and methods**

108 The major innovations in PING, detailed below, include the use of multiple-sequence per gene
109 alignment references that incorporate the allelic diversity of *KIR*, and genotype-matched alignment
110 references (Figure 1). The use of a diverse reference set in the copy number module substantially
111 improves copy number determination for *KIR2DL2*, *KIR2DL3*, *KIR2DS1* and *KIR2DL1* compared
112 to our prototype approach. The improved performance enables a supplementary alignment
113 workflow that dynamically constructs genotype-matched alignment references based on the so-
114 established gene content and a preliminary genotype determination. The genotypes determined by
115 this novel genotype-aware alignment workflow are highly accurate, with few unresolved calls.

116



117

118 **Figure 1. Overview of the PING pipeline.** The PING pipeline processes *KIR* targeted
119 sequencing data to determine *KIR* gene copy number and allele genotypes through a series of
120 modules. First, *KIR* aligned reads are filtered through an alignment to a set of *KIR* haplotypes in
121 PING extractor. Second, copy number of *KIR* genes are determined through an exhaustive
122 alignment to a diverse set of *KIR* sequences in PING copy. Finally, PING allele performs a series

123 of alignments to determine the most congruent *KIR* genotype, which informs a final round of
124 alignment and genotype determination. Additionally, PING reports any identified novel SNPs
125 and new alleles (SNP combinations not found in any described *KIR* allele sequence).

126

127 K-mer similarity analysis

128 Read misalignments due to sequence identity across the *KIR* region are a persistent challenge in
129 *KIR* bioinformatics and often lead to spurious genotyping results. To quantify the extent of
130 sequence identity and inform our investigation of SNPs suspected to be originating from
131 misaligned reads, we performed a *k*-mer similarity analysis using all 905 described *KIR* allele
132 sequences in the Immuno Polymorphism Database (IPD) - *KIR* (21), release 2.7.1. Here, we
133 transformed allele sequences into all distinct subsequences of length *k* to compare sequence
134 identity between genes.

135

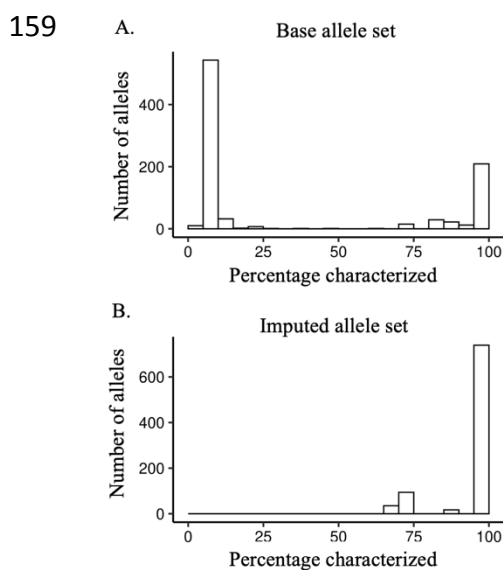
136 Considering all allele sequences of a given *KIR* gene, we generated all distinct k-mers of sizes 50,
137 150 and 250 nucleotides and counted the number of exact k-mer matches across *KIR* genes for
138 each value of *k*. Shared k-mer proportions were calculated by dividing the number of shared k-
139 mers by the total number of k-mers of that gene. K-mer sharing diagrams were generated using
140 the circlize (34) package in R (35). Since this approach displays connections between two genes,
141 k-mers that were shared between three or more *KIR* genes were enumerated by dual gene
142 combinations, so that a k-mer that matches genes A, B and C count as matches between A-B, A-
143 C and B-C.

144

145 Imputation of uncharacterized regions and extension of untranslated regions to generate
146 comprehensive alignment reference sequence

147 Sequence alignment workflows that fail to account for the allelic diversity of *KIR*, such as
148 workflows using single-sequence per gene references, increase the likelihood that read alignments
149 will be biased towards an incorrect gene. Conversely, incorporating the allelic diversity of *KIR*
150 into the alignment reference enables the identification and processing of reads that map to multiple
151 *KIR* genes. However, because the sequences of many *KIR* alleles have only been characterized for
152 exons, 65% of named alleles have less than 20% of their full-length sequence characterized (Figure
153 2A). Additionally, IPD-KIR allele sequences only include ~250bp of 5' untranslated region (UTR)
154 sequence and ~500bp of 3' UTR sequence, reducing alignment depths across the first exon and
155 potential regulatory regions. Thus, to maximize the utility of reference *KIR* sequences, we
156 designed and implemented a protocol to impute the sequence of the intronic regions, followed by
157 a protocol to extend UTRs to 1000bp each.

158



160

161 **Figure 2. *KIR* sequence characterization before and after imputation.** (A) Histogram of IPD-
162 *KIR* allele sequence lengths, shown as percentage of longest sequence for each gene and major
163 allele group. (B) Histogram of IPD-*KIR* allele sequence lengths after imputation, shown as
164 percentage of longest sequence for each gene and major allelic group.

165

166 As reference sequences, we used all *KIR* alleles described in the IPD - *KIR* (21), release 2.7.1. A
167 subset of these sequences is not completely characterized through all exons and introns. We
168 therefore used gene-specific alignments of known sequences, provided by IPD-*KIR* as multiple
169 sequence format (MSF) files, and completed each allele sequence to comprise the invariant
170 nucleotides together with each variable position represented by an 'N'. Using this imputation
171 method, we generated a new set of reference alleles in which ~90% of the 905 alleles were >98%
172 complete (Figure 2B).

173

174 To extend UTR sequence, donor sequences, sourced from the full *KIR* haplotype sequences used
175 in PING extractor, were appended to the ends of each reference sequence to generate 1000bp long
176 UTRs. A single 3'UTR and 5'UTR donor sequence was used for each gene and major allelic group.

177

178 Copy number determination workflow

179 The high sequence similarity between *KIR* genes coupled with extensive structural variation and
180 nucleotide diversity makes copy number determination a non-trivial task. Our copy determination
181 method is largely identical to that described in Norman et al. (27), in which copy number is
182 determined by comparing the number of reads that align uniquely to each *KIR* gene across a batch
183 of samples using *KIR3DL3* as a normalizer. The improvement made by our method is the use of a

184 comprehensive *KIR* reference composed of 905 distinct sequences from the imputed and extended
185 allele set, instead of a single-sequence per gene reference. The use of a comprehensive reference
186 provides a more accurate comparison of the number of reads that align uniquely to any *KIR* gene.

187

188 *KIR* virtual probes

189 To determine the presence of target alleles or allelic groups that are prone to misidentification due
190 to read misalignments, we have developed a set of virtual, or text-based, probes. The probe set
191 includes those described in Norman et al. (27), as well as additional, custom probes (Table in S5
192 Table). Probes are designed to match sequence that is unique to the target allele or allelic group,
193 and sequence uniqueness is determined by a grep search over the imputed and extended IPD-KIR
194 sequence set. Application of the probe set is performed using grep over the sequencing data,
195 counting the number of unique reads that contain sequence perfectly matching the probe. A probe
196 hit is determined using a threshold of 10 matching reads.

197

198 Designing a minimized reference allele set

199 While performing a comprehensive alignment to the full *KIR* allele set reduces misalignments
200 caused by reference sequence bias, it demands substantial resource utilization, as well as a large
201 alignment and processing time cost which can prove untenable for processing large datasets. For
202 example, copy determination processing for 10 paired-end sequences using 36 threads took 4.95
203 hours with a maximum Binary Alignment Map (BAM) (36) file size of 338.2MB.

204

205 To address this issue, we constructed a minimized set of reference alleles to improve resource
206 utilization, alignment and processing times while still reducing misalignments caused by reference

207 sequence bias. The minimized reference set consists of five alleles for each *KIR* gene and major
208 allelic group (Table in S7 table). The use of five alleles per gene was empirically determined to be
209 sufficient for reducing reference sequence bias, while still considerably reducing the
210 computational burden of multiple-sequence per gene alignments. Designing this reference set was
211 guided by selecting alleles which had fully-characterized or nearly fully-characterized sequence,
212 the secondary criteria was maximizing SNP diversity between the reference alleles of each gene,
213 and third was selecting reference alleles to sequester reads susceptible to off-gene mapping. For
214 example, reference sequences to represent *KIR2DS1*002* as well as *KIR2DL1*004* were selected
215 to sequester reads that perfectly align to both. Notable characteristics of the reference set are the
216 separation of *KIR2DL2* from *KIR2DL3*, the separation of *KIR3DL1* from *KIR3DS1*, and the
217 merging of *KIR2DL5A* and *KIR2DL5B*.

218

219 Genotype determination workflow

220 Indexed reads, detailed in S1 text, are processed to generate a depth table spanning -1000bp 5'UTR
221 to 1000bp 3'UTR for each *KIR* gene and major allelic group. Depths are marked independently for
222 A, T, C, G, deletions and insertions. Depth tables are processed to generate SNP tables for positions
223 passing a minimum depth threshold (default 8 for initial genotyping and 20 for final genotyping).
224 To identify heterozygous positions the depth of each aligned variant is divided by the highest depth
225 variant for that position, and up to three variants (A, T, C, G, deletions and insertions) passing the
226 ratio threshold (default 0.25 for initial and final genotyping) are recorded.

227

228 Genotypes for each gene and major allelic group are determined from the aligned SNPs using a
229 mismatch scoring approach. First, aligned homozygous SNPs are compared to each IPD-KIR

230 allele, with SNP mismatches counting as a score of 1 and matches as 0. The lowest scoring alleles
231 and alleles within a set scoring buffer of the lowest score (default of 4 for the initial genotyping
232 workflow and 1 for the final genotyping workflow), are carried over into heterozygous position
233 scoring.

234

235 For aligned heterozygous position scoring, all possible allele combinations are enumerated
236 according to the determined copy of the gene under consideration, up to copy 3. For each aligned
237 position, the variant(s) for each allele combination are compared to the aligned variants, with full
238 matches counted as a score of 0 and mismatches scored according to the number of mismatched
239 variants. For each allele combination, the homozygous score of each component allele is added to
240 the heterozygous score, and the lowest scoring combinations are returned as the determined
241 genotype. For the final genotyping workflow, only perfectly scoring combinations are accepted,
242 with any mismatches resulting in an unresolved genotype.

243

244 The same workflow is applied to both initial and final genotyping with some important
245 distinctions. In the initial genotyping workflow, the imputed and extended IPD-KIR allele
246 sequences are used for SNP comparisons, uncharacterized variants within the comparison
247 sequences are marked as full mismatches, and all aligned allele-differentiating SNP positions
248 passing the depth threshold are compared and used for scoring. In the final genotyping workflow,
249 the unimputed IPD-KIR allele sequences are used for SNP comparisons, uncharacterized variants
250 within the comparison sequences are marked as matches, and only aligned exonic SNP positions
251 passing the depth threshold are compared and used for scoring.

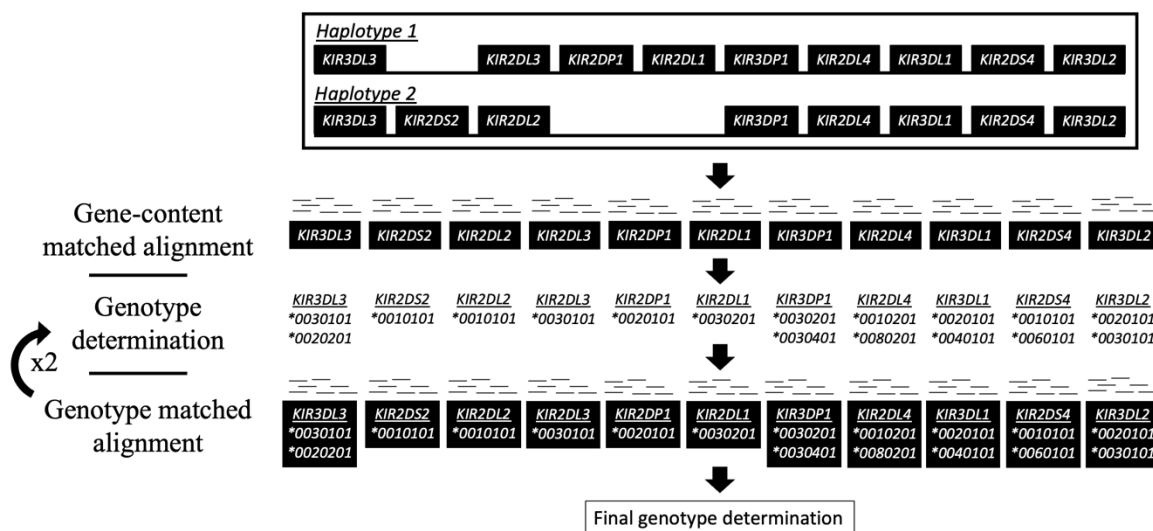
252

253 The final exonic resolution genotypes are processed to add null alleles to the genotype string for
 254 genes with copy 0 or copy 1, and combine component allele typings for the major allelic groups
 255 *KIR2DL2* and *KIR2DL3*, *KIR3DL1* and *KIR3DS1*, and *KIR2DS3* and *KIR2DS5*.

256

257 Genotype matched alignment workflow

258 The overall alignment strategy of PING is to reduce reference sequence bias through the use of
 259 multiple-sequence per gene references, and the use of references that reflect the gene content
 260 makeup or genotype makeup of a sample (Figure 3). Additionally, PING utilizes multiple rounds
 261 of alignment and genotype determination with varied processing parameters to reduce bias
 262 introduced by assumptions made during the processing workflow.



263

264 **Figure 3. Overview of the genotype aware alignment workflow.** Sequence gene content,
 265 determined by PING copy, informs the selection of reference sequence from a predefined set of
 266 diverse allele sequences. An exhaustive alignment is performed to the selected allele set, from
 267 which an initial genotype determination is made. The determined genotype informs selection of
 268 reference alleles for a genotype aware alignment, followed by another round of genotype
 269 determination. The genotype aware alignment and subsequent genotype determination is repeated,

270 and the most congruent genotypings across all alignment rounds inform reference selection for a
271 final round of alignment. A non-exhaustive alignment is performed to the selected allele set, from
272 which all aligned reads are processed and used for the final genotype determination.

273

274 The first step is an alignment to a multiple-sequence per gene, gene content matched reference.
275 This gene-content aware alignment workflow constructs individualized alignment references
276 based on the presence of certain *KIR* genes: *KIR3DP1*, *KIR2DS2*, *KIR2DL23*, *KIR2DL5A*,
277 *KIR2DL5B*, *KIR2DS3*, *KIR2DS5*, *KIR2DP1*, *KIR2DL1*, *KIR2DL5*, *KIR3DL1S1*, *KIR2DS4*,
278 *KIR3DL2*, *KIR2DS1* and *KIR3DL3* (assumed always present (37)). Reference sequences are
279 selected from the diverse, minimized reference sequence set described above. We have included
280 an option to align to the full comprehensive allele set, but this is not default behavior as these
281 alignments are time and resource intensive.

282

283 An exhaustive alignment, an alignment in which all qualified read mappings are recorded, is
284 performed and aligned reads are processed and formatted according to the alignment processing
285 workflow, detailed in S1 text, selecting for reads that uniquely map to a gene or major allelic
286 group. The formatted uniquely-mapped read set is processed according to the genotype
287 determination workflow to obtain an initial full-resolution genotype determination.

288

289 Second is a series of two alignments to genotype-matched references with varied processing
290 parameters to identify the most congruent *KIR* genotype. Genotype congruence is determined by
291 the least number of SNP mismatches between the determined allele typing(s) of a gene, and the
292 aligned SNPs. For each genotype-matched alignment in this series, the determined allele typing(s),

293 including any ambiguity, are used as reference sequence for the following alignment. For each
294 alignment, genotypes are first determined at seven-digit (non-coding mutation level), then five-
295 digit resolution (synonymous mutation level). This approach reduces the impact of uncharacterized
296 regions of IPD-KIR allele sequences on genotype determination, as most sequences are fully
297 characterized across exons. Genotype determination can be biased towards or against IPD-KIR
298 alleles with uncharacterized regions depending on whether uncharacterized SNPs count as
299 mismatches or not. To reduce time spent on genotype determination, any unambiguous typing that
300 is perfectly matched to the aligned SNPs is locked in across all subsequent intermediate rounds of
301 genotype determination.

302

303 The reference for the final alignment is built from the locked genotypings and the closest matched
304 genotypings for genes without a locked genotyping. A non-exhaustive alignment is performed to
305 the built reference, from which all aligned reads are processed and formatted according to the
306 alignment processing workflow. The formatted read alignments are passed to the genotype
307 determination workflow to obtain a final exonic (five-digit) resolution genotype determination.

308

309 Additional methods utilized by the genotype matched alignment workflow

310 Genotype matched alignments and subsequent genotype determinations can get stuck on a
311 mistyped allele due to persistent reference sequence bias. In other words, a false SNP call
312 originating from misaligned reads can perpetuate itself in the genotype matched alignments due to
313 the same allele determination being made and the same alignment reference being used. To address
314 this issue, we have included a method in the genotype matched alignments that will add the five
315 allele sequences from the diverse, minimized reference set to the genotype-matched reference for

316 any gene with an allele typing that does not perfectly match the aligned SNPs. The rationale behind
317 this method is that mismatched allele typings are likely due to misaligned reads, and the use of the
318 mismatched allele sequence as a reference will cause the read misalignments to be repeated in
319 subsequent alignment and genotyping rounds. The addition of a diverse set of alignment sequences
320 gives an avenue to break from this cycle by increasing the likelihood that a different allele typing
321 will be made.

322

323 In building genotype-matched references PING allows any allele to be used as reference sequence,
324 however, some allele sequences are only partially characterized even after imputation. The use of
325 allele sequences containing uncharacterized sequence as alignment references can introduce
326 reference sequence bias and drive read misalignments even if the reference alleles perfectly match
327 the true genotype of the sample. To address this issue, we have included a method to add fully-
328 characterized sequence to the alignment reference for any gene represented by only partially-
329 characterized sequence(s). Fully-characterized alleles are pulled from the diverse, minimized
330 reference set.

331

332 In the genotype-aware alignment workflow we found issues with false negative identifications of
333 *KIR2DL1*004/*007/*010* due to reads cross-mapping to other gene sequences. This issue was
334 rectified using virtual sequence probes specific to each of these *KIR2DL1* allele groups to identify
335 **004/*007/*010* allele presence. If *KIR2DL1*010* is present, then the *KIR2DL1*010* allele
336 sequence is added to the alignment reference. If *KIR2DL1*004* is present, then the
337 *KIR2DL1*0040101* allele sequence is added to the alignment reference. If *KIR2DL1*007* is
338 present, then the *KIR2DL1*007* allele sequence is added to the alignment reference. If multiple of

339 these allele groups are present, then the *KIR2DL1*0040101* allele sequence is added to the
340 alignment reference.

341

342 We implemented additional probes to identify alleles and structural variants prone to
343 misidentification across *KIR2DL1*, *KIR2DL2*, *KIR2DL4*, *KIR2DS1*, *KIR3DP1* and *KIR3DS1*. For
344 example, we implemented a probe to identify the *KIR2DL4* poly-A stretch at the end of exon 7, as
345 well as a probe to identify *KIR3DP1* exon 2 deletion variants. The full list of probes used for
346 reference refinement can be found in S5 table.

347

348 *KIR* synthetic sequence dataset

349 A *KIR* synthetic dataset consisting of 50 sequences was generated using the ART next-generation
350 sequencing read simulator (38). ART parameters were set to simulate 150-bp paired-end reads at
351 50x coverage, with a median DNA fragment length of 200 using quality score profiles from the
352 HiSeq 2500 system. Eleven of the *KIR* haplotypes described in Jiang et al. (39) were used to
353 simulate structural variation of the *KIR* region. Two of the eleven haplotypes were randomly
354 selected with replacement to establish the copy number for each sample. Allele sequences were
355 selected randomly without replacement from the imputed and extended set according to the copy
356 number of each gene. Any uncharacterized regions in the selected allele sequences were replaced
357 with sequence from a random fully-characterized sequence from the same gene. Reads were named
358 according to the source allele, enabling tracing of misaligned reads to their source allele and gene.
359 The full synthetic dataset is available at: https://github.com/wesleymarin/KIR_synthetic_data.

360

361 Discordant genotype results for the synthetic dataset were investigated by identifying the source
362 gene for each read aligned to the incorrectly genotyped gene. The results were summarized to show
363 the total number of reads from each source gene to each aligned gene, Table A in S8 Table, and a
364 read sharing diagram was generated using the circlize (34) package in R (35).

365

366 Characterization of *KIR* reference cohorts for PING development

367 A significant barrier to the development of bioinformatic methods for high-resolution *KIR*
368 sequence interpretation is the lack of a well-characterized reference cohort. Without such a
369 resource it is extremely difficult to recognize and resolve issues with read misalignments, which
370 can result in SNP calls that appear reasonable in many cases. To resolve this issue, we have
371 characterized a *KIR* reference cohort of 379 healthy individuals of European ancestry that had been
372 previously sequenced using our *KIR* target capture method (40), with the results meticulously
373 curated by manual alignment and inspection of all sequences to provide a ground truth dataset to
374 aid pipeline development (Table A in S4 Table). Furthermore, the European samples were
375 independently sequenced and genotyped for *KIR* by our collaborators at the DKMS registry for
376 volunteer bone marrow donors (31). Any discordant typing or gene content results were resolved
377 through direct examination of sequence alignments, and where necessary, confirmatory
378 sequencing.

379

380 In order to validate our method on a second, divergent population, we also examined a previously
381 characterized cohort of African Khoisan individuals (41), for which *KIR* alignments and genotypes
382 were manually inspected (Table B in S4 Table).

383

384 Copy number and genotype concordance calculations

385 For the European cohort, any genotype containing an unresolved *KIR3DL3* genotyping, a genotype
386 for which the aligned SNPs do not perfectly match currently described alleles, in the truth dataset
387 were excluded from copy number and genotype concordance comparisons. There were 16 full
388 genotypes excluded by this criterion. Additionally, any individual gene with an unresolved
389 genotype in the truth dataset were excluded from copy number and genotype concordance
390 comparisons. For the synthetic dataset there were no simulated new alleles, so the full dataset was
391 used for copy number and genotype concordance comparisons. For the Khoisan dataset any
392 individual gene with an unresolved genotype in the truth dataset were excluded from genotype
393 concordance comparisons but were included for copy number comparisons.

394

395 Copy concordance was calculated by directly comparing the determined copy values to the
396 validation copy values (Tables in S2 Table). Genotype concordance was calculated on a per-gene
397 basis by comparing each component allele of the determined typing to the truth genotype.

398

399 Code availability

400 The PING pipeline is available at <https://github.com/wesleymarin/PING> (42) with the following
401 open source license: <https://github.com/wesleymarin/PING/blob/master/LICENSE>. Scripts and
402 datasets used for data analysis are available at https://github.com/wesleymarin/ping_paper_scripts.
403 PING was developed in the R programming language and tested on a Linux system. Additional
404 requirements are: Samtools v1.7 or higher (36), Bcftools v1.7 or higher, and Bowtie2 v2.3.4.1 or
405 higher (43). The synthetic KIR sequence dataset is available at
406 https://github.com/wesleymarin/KIR_synthetic_data.

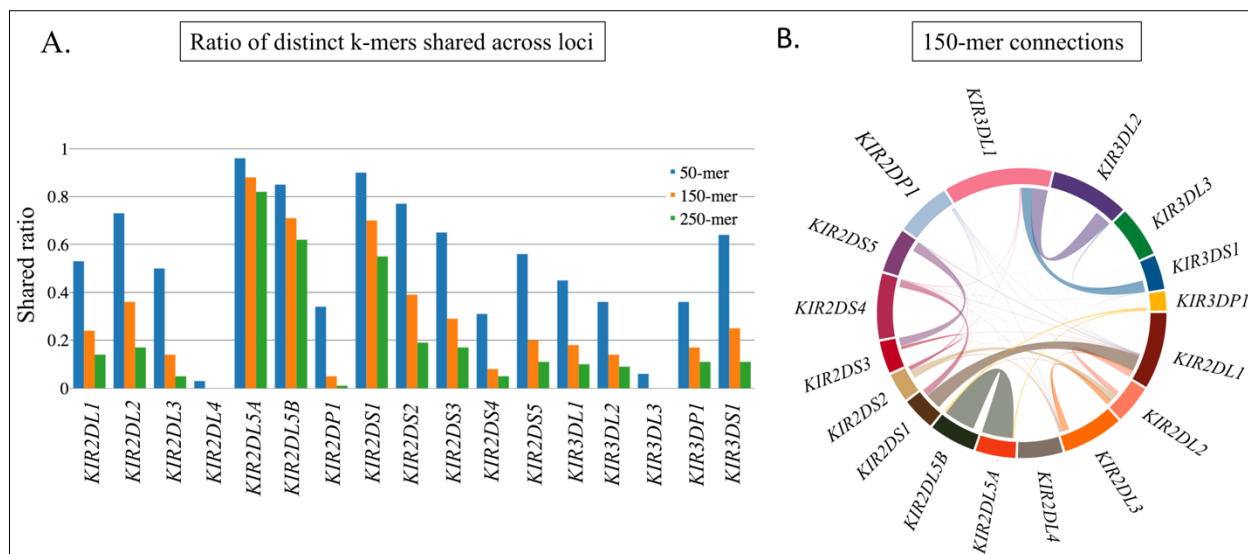
407

408 Results

409 Extensive sequence identity among *KIR* genes is a major barrier for interpreting short-read
410 sequencing data

411 Our analysis showed that many genes share significant sequence identity at sequencing lengths
412 commonly used in next generation sequencing (NGS) technology (Figure 4A). For example,
413 *KIR2DL5A* shares 12,591 of its 15,359 distinct 150-mers (82%) with *KIR2DL5B* (Figure 4B,
414 Tables in S6 Table), making it extremely difficult to distinguish NGS (next generation sequencing)
415 short reads originating from these genes. Likewise, over 90% of the distinct 50-mers, and over
416 50% of distinct 250-mers of *KIR2DS1* are shared with other genes, the vast majority with
417 *KIR2DL1*. This analysis allowed us to identify specific “hotspots” for read misalignments,
418 informing post-alignment modifications (detailed previously) to minimize their impact.

419



420

421 **Figure 4. K-mer analysis of *KIR* gene sequence similarity.** (A) Ratio of distinct k-mers of size
422 50, 150 and 250 that are shared between the indicated *KIR* gene and others. The inverse of these

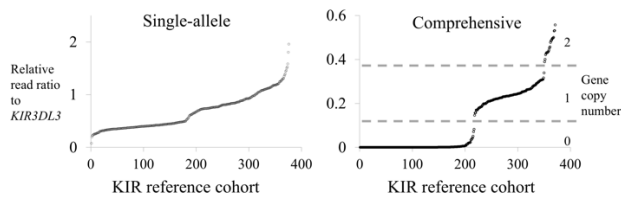
423 bars (not shown) would indicate the proportion of k-mers that are distinct to that gene and not
424 found in the alleles of other genes. (B) 150-mer connections between *KIR* genes, the size of the
425 connecting line roughly indicates the total number of shared 150-mers.

426

427 Development of a comprehensive *KIR* alignment reference enables accurate copy number
428 determination of *KIR2DL1*, *KIR2DS1*, *KIR2DL2* and *KIR2DL3*.

429 Applying the comprehensive reference allele set to gene content and copy number determination,
430 we achieved significant improvement over a single-sequence per gene reference for *KIR2DL1*,
431 *KIR2DS1* and the allelic groups *KIR2DL2* and *KIR2DL3* (S1 Figure, S2 Figure and S3 Figure).
432 The copy number for *KIR2DS1*, per example, which is highly prone to read misalignments due to
433 similarity to *KIR2DL1* and *KIR2DS4* (Figure 4B), was clearly determined (Figure 5).

434



435

436 **Figure 5. Use of comprehensive reference can improve copy determinations.** Single-
437 sequence reference vs. comprehensive reference copy number plot of *KIR2DS1*. The copy plot of
438 the single-sequence reference alignment shows no differentiation between copy groupings while
439 the comprehensive reference alignment shows a clear distinction between the copy 0, 1 and 2
440 groups.

441

442 PING delivers accurate copy number and high-resolution allele calls

443 The overall performance of PING was assessed using our European *KIR* reference cohort, a
444 synthetic *KIR* dataset, and a Khoisan *KIR* reference cohort. Results for PING copy number
445 determination are summarized in Table 1, showing at least 97% concordance for the European
446 cohort for all compared genes, with most genes exhibiting more than 99% concordance.
447 Performance for the synthetic dataset showed 100% copy concordance for all compared genes
448 except for *KIR2DL1*, at 98%, and *KIR2DS3*, at 92%. Finally, performance for the Khoisan cohort
449 showed at least 95% concordance for all compared genes except for *KIR2DL2*, at 61%,
450 *KIR2DL5A/B*, at 88%, *KIR2DS5*, at 89%, and *KIR2DL1*, at 94%. Across all datasets *KIR3DL3* was
451 not compared due to its use as a reference gene, and for the European and Khoisan cohorts the
452 pseudogene *KIR3DPI* was not compared due to an absence of validation data.

453

<u>Gene</u>	<u>European</u>	<u>N</u>	<u>Synthetic</u>	<u>N</u>	<u>Khoisan</u>	<u>N</u>
<i>KIR3DL3</i>	-	-	-	-	-	-
<i>KIR2DS2</i>	0.988	343	1.00	50	0.97	100
<i>KIR2DL2</i>	0.994	331	1.00	50	0.61	100
<i>KIR2DL3</i>	0.994	331	1.00	50	1.00	100
<i>KIR2DL5A/B</i>	0.997	343	1.00	50	0.88	100
<i>KIR2DS3</i>	0.988	343	0.92	50	0.97	100
<i>KIR2DS5</i>	0.985	342	1.00	50	0.89	100
<i>KIR2DPI</i>	0.982	338	1.00	50	0.95	100
<i>KIR2DL1</i>	0.970	334	0.98	50	0.94	100
<i>KIR3DPI</i>	-	-	1.00	50	-	-
<i>KIR2DL4</i>	0.994	341	1.00	50	1.00	100
<i>KIR3DL1</i>	0.997	340	1.00	50	1.00	100
<i>KIR3DS1</i>	0.997	339	1.00	50	0.99	100
<i>KIR2DS1</i>	0.988	342	1.00	50	1.00	100
<i>KIR2DS4</i>	0.991	343	1.00	50	0.99	100
<i>KIR3DL2</i>	0.988	326	1.00	50	1.00	100

454

455 **Table 1. Copy number determination performance.** Concordance table comparing copy
 456 numbers determined by PING for the European reference cohort, a synthetic *KIR* dataset, and a
 457 Khoisan reference cohort.

458
 459 Performance of genotype determination was assessed at three-digit resolution (protein level) for
 460 the European and Khoisan cohorts, and at five-digit resolution (synonymous mutation level) for
 461 the synthetic dataset (Table 2). The results were categorized as genotype matches, mismatches or
 462 unresolved genotypes, which were cases where PING could not make a genotype determination.

463
 464 PING genotype determination for the European cohort showed low percentages of unresolved
 465 genotypes with few mismatches for all compared genes except for *KIR2DP1*, with 10.3%
 466 unresolved. Notable results for the European cohort were the low frequencies of unresolved
 467 genotypes across most genes, except *KIR2DP1*, and the extremely low frequencies of mismatched
 468 genotypes, below 1%, for 9 out of the 12 genes compared.

Gene	Dataset	Match	Mismatch	Unresolved	N
<i>KIR3DL3</i>	European	0.959	0.009	0.032	686
	Synthetic	0.960	0.000	0.040	100
	Khoisan	0.887	0.062	0.050	80
<i>KIR2DS2</i>	European	0.975	0.012	0.013	686
	Synthetic	0.940	0.040	0.020	100
	Khoisan	0.835	0.005	0.160	188
<i>KIR2DL23</i>	European	0.965	0.003	0.032	656
	Synthetic	0.850	0.020	0.130	100
	Khoisan	0.810	0.042	0.149	168
<i>KIR2DL5A/B</i>	European	0.927	0.044	0.029	687
	Synthetic	0.856	0.106	0.038	104
	Khoisan	0.857	0.071	0.071	126
<i>KIR2DS35</i>	European	0.982	0.006	0.012	683
	Synthetic	0.923	0.019	0.058	104
	Khoisan	0.871	0.052	0.078	116
<i>KIR2DP1</i>	European	0.890	0.007	0.103	672
	Synthetic	0.971	0.000	0.029	103
	Khoisan	0.721	0.012	0.267	86
<i>KIR2DL1</i>	European	0.961	0.009	0.030	666
	Synthetic	0.883	0.019	0.097	103

	Khoisan	0.803	0.045	0.152	132
<i>KIR3DP1</i>	European	-	-	-	-
	Synthetic	0.773	0.055	0.173	110
	Khoisan	-	-	-	-
<i>KIR2DL4</i>	European	0.980	0.007	0.013	685
	Synthetic	1.000	0.000	0.000	110
	Khoisan	0.961	0.006	0.032	154
<i>KIR3DL1S1</i>	European	0.962	0.006	0.032	686
	Synthetic	0.955	0.000	0.045	110
	Khoisan	0.873	0.028	0.099	142
<i>KIR2DS1</i>	European	0.985	0.003	0.012	682
	Synthetic	0.900	0.000	0.100	100
	Khoisan	0.995	0.000	0.005	198
<i>KIR2DS4</i>	European	0.943	0.044	0.013	685
	Synthetic	0.940	0.020	0.040	100
	Khoisan	0.793	0.051	0.157	198
<i>KIR3DL2</i>	European	0.965	0.008	0.028	648
	Synthetic	0.990	0.010	0.000	100
	Khoisan	0.819	0.011	0.170	188

469

470 **Table 2. Genotype determination performance.** Genotype determination performance table

471 comparing the genotypes determined by PING to the validation genotypes for each dataset.

472 Possible outcomes are ‘Match’, where the determined component allele matches the validation

473 allele, ‘Mismatch’, where the determined component allele does not match the validation allele, or

474 ‘Unresolved’, where PING was unable to determine a genotype, but the validation allele was not

475 marked as unresolved. The coloring signifies concordance level, where green is 0-10% discordant,

476 yellow is 10-15% discordant, and red is over 15% discordant.

477

478 Determined genotypes for the synthetic dataset showed over 95% concordance for *KIR3DL3*,

479 *KIR2DP1*, *KIR2DL4*, *KIR3DL1S1*, and *KIR3DL2*. However, the synthetic dataset showed high

480 percentages of unresolved genotypes for *KIR2DL23*, *KIR3DP1* and *KIR2DS1*, each over 10%

481 unresolved, and *KIR2DS35* and *KIR2DL1* showed 5.8% and 9.7% unresolved, respectively.

482 *KIR2DL5A/B* and *KIR3DP1* showed the highest mismatched genotype percentages, at 10.6% and

483 5.5%, respectively, while *KIRDL3*, *KIR2DP1*, *KIR2DL4* and *KIR2DS1* each showed 0.0%
484 mismatched genotypes.

485

486 Determined genotypes for the Khoisan cohort showed highly concordant genotypes for *KIR2DL4*,
487 at 96.1%, and *KIR2DS1*, at 99.5%. Additionally, results for this dataset showed low mismatch
488 frequencies for *KIR2DS2*, *KIR2DL23*, *KIR2DP1*, *KIR2DL1*, *KIR3DL1S1* and *KIR3DL2*, each
489 below 5.0% mismatched. However, the Khoisan cohort showed moderate mismatch frequencies
490 for *KIR3DL3*, at 6.2%, *KIR2DL5A/B*, at 7.1%, *KIR2DS35*, at 5.2%, and *KIR2DS4*, at 5.1%, and
491 higher unresolved rates for *KIR2DS2*, *KIR2DL23*, *KIR2DP1*, *KIR2DL1*, *KIR2DS4* and *KIR3DL2*,
492 each over 10.0% unresolved.

493

494 For the European and Khoisan cohorts the pseudogene *KIR3DP1* was not compared due to an
495 absence of validation data.

496

497 Looking specifically at the concordance of resolved genotypes, the European cohort showed
498 greater than 98.0% concordance across all compared genes except for *KIR2DL5A/B*, at 95.5%, and
499 *KIR2DS4*, at 95.6% (Table 3). The synthetic dataset showed 100% concordance for *KIR3DL3*,
500 *KIR2DP1*, *KIR2DL4*, *KIR3DL1S1* and *KIR2DS1*, over 95% concordance for *KIR2DS2*,
501 *KIR2DL23*, *KIR2DS35*, *KIR2DL1*, *KIR2DS4* and *KIR3DL2*. The lowest performing genes in the
502 synthetic dataset were *KIR2DL5A/B*, at 89%, and *KIR3DP1*, at 93%. The Khoisan cohort showed
503 100% concordance for *KIR2DS1*, over 95% concordance for *KIR2DS2*, *KIR2DL23*, *KIR2DP1*,
504 *KIR2DL4*, *KIR3DL1S1* and *KIR3DL2*, and over 90% concordance for *KIR3DL3*, *KIR2DL5A/B*,
505 *KIR2DS35*, *KIR2DL1* and *KIR2DS4*.

506

<u>Gene</u>	<u>European</u>	<u>N</u>	<u>Synthetic</u>	<u>N</u>	<u>Khoisan</u>	<u>N</u>
<i>KIR3DL3</i>	0.991	664	1.00	96	0.934	76
<i>KIR2DS2</i>	0.988	677	0.96	98	0.994	158
<i>KIR2DL23</i>	0.997	635	0.98	87	0.951	143
<i>KIR2DL5A/B</i>	0.955	667	0.89	100	0.923	117
<i>KIR2DS35</i>	0.994	675	0.98	98	0.944	107
<i>KIR2DP1</i>	0.992	603	1.00	100	0.984	63
<i>KIR2DL1</i>	0.991	646	0.98	93	0.946	112
<i>KIR3DP1</i>	-	-	0.93	91	-	-
<i>KIR2DL4</i>	0.993	676	1.00	110	0.993	149
<i>KIR3DL1S1</i>	0.994	664	1.00	105	0.969	128
<i>KIR2DS1</i>	0.997	674	1.00	90	1.000	197
<i>KIR2DS4</i>	0.956	676	0.98	96	0.940	167
<i>KIR3DL2</i>	0.992	630	0.99	100	0.987	156

507

508 **Table 3. Resolved genotype concordance.** PING genotype determination performance for the
509 European reference cohort, a synthetic *KIR* dataset, and the Khoisan reference cohort for each
510 considered *KIR* gene.

511

512 Together, these results demonstrate that PING accurately provides *KIR* genotyping across distinct
513 populations.

514

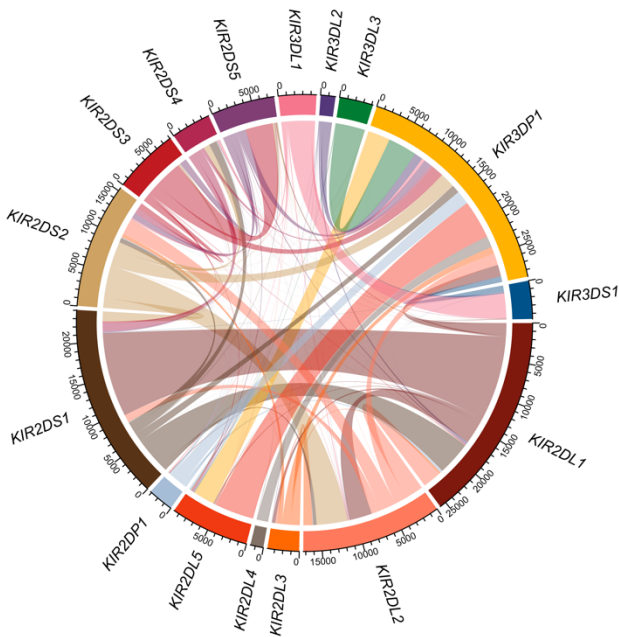
515 Analysis of discordant determined copy number and genotype results

516 The discordant copy results for *KIR2DS3* in the synthetic dataset were the result of poor
517 differentiation between copy groups (S3 Figure). The highly discordant *KIR2DL2* copy number
518 result for the Khoisan cohort was due to non-differentiable copy number groupings (S2 Figure).
519 Since the *KIR2DL3* copy differentiation for this cohort was well defined, these results were used

520 to set the *KIR2DL2* copy number prior to genotype determination using the formula
521 $KIR2DL2_copy = 2 - KIR2DL3_copy$.

522

523 An investigation into the discordant genotypes for the synthetic dataset showed discordant
524 genotype determination results for *KIR2DS35* were largely due to source reads from *KIR2DS3*
525 aligning to *KIR2DS5* reference sequence, with a smaller number of reads from *KIR2DS5* aligning
526 to *KIR2DS3* reference sequence (Figure 6, Table A in S8 Table). This differential read flow
527 between the two allelic groups is reflected in the component allele typings, with six discordant
528 *KIR2DS5* genotypings and two discordant *KIR2DS3* genotypings (Table C in S3 Table). Intragenic
529 misalignments are a product of how the PING workflow is structured, as major allelic groups, such
530 as *KIR2DS3* and *KIR2DS5*, are treated as independent genes during alignment and genotyping.
531 Intragenic misalignments were also a large contributor to *KIR3DL1S1* discordance, with reads
532 supplied by *KIR3DL1* mapping to *KIR3DS1* reference sequence.



533

534 **Figure 6. Misaligned read sources in the synthetic dataset.** Analysis of mismatched or
535 unresolved genotype determination results for the synthetic sequence dataset where all
536 misaligned reads are traced back to their source gene. The connections between genes represent
537 the number of misaligned reads, and the color of the connections represent the source gene.

538

539 The analysis showed *KIR3DP1* as a major hub for receiving misaligned reads, with reads being
540 contributed by each other *KIR* gene. In fact, *KIR3DP1* was largely the only receiver for misaligned
541 reads originating from *KIR3DL3*, *KIR3DL2*, *KIR2DP1* and *KIR2DL4*. While the only genes
542 receiving reads sourced from *KIR3DP1* were *KIR2DL5A* and *KIR2DL5B*.

543

544 The analysis also showed several gene pairings, where two genes largely sent and received reads
545 from one another. Once such pairing was between *KIR2DL1* and *KIR2DS1*, where each gene were
546 the largest contributor and receiver of reads for each other. Another pairing was between *KIR2DL2*
547 and *KIR2DS2*, although both genes sent and received reads from several other genes.

548

549 This analysis illustrates the complex and highly interconnected nature of *KIR* and highlights the
550 difficulty behind accurate interpretation of *KIR* short-read sequencing data.

551

552 Performance

553 The run time and resource utilization of the PING pipeline was measured on an Intel Xeon 2.20
554 GHz CPU using 36 threads. For ten sequences from the synthetic dataset, it took 1.92 mins for
555 *KIR* read extraction, 34.7 mins for copy determination aligning to the minimized reference set, and

556 2.10 hours for genotype determination aligning to the minimized reference set. The output
557 directory size was 1.4GB.

558

559 **Discussion**

560 Our k-mer analysis of all documented *KIR* variation shows the high degree of sequence identity
561 between *KIR* genes and illustrates the challenges imposed by the homology of *KIR* on short-read
562 interpretation workflows. It demonstrates that some genes are more likely to exhibit read
563 misalignment problems than others. *KIR2DP1*, *KIR3DL3* and *KIR2DL4* have relatively unique
564 sequence, while *KIR2DS1*, *KIR2DL5A* and *KIR2DL5B* have considerable shared sequence. This
565 type of analysis provides an informative tool for investigating irregularities in the processing of
566 *KIR* sequence data, revealing which genes are likely to be erroneously interpreted due to read
567 misalignments for common sequencing read lengths. While paired-end sequencing with longer
568 reads can improve read alignment fidelity, in our own experience 290bp paired-end reads with a
569 median insert length of approximately 600bp still exhibited considerable read misalignment
570 problems. It is important to note that this analysis does not account for unknown variation or
571 intergenic sequence, two other sources of sequence variation that could potentially result in
572 misaligned reads.

573

574 An initial determination of *KIR* gene content and copy number provides an informative scaffold
575 for minimizing misalignments through the exclusion of reference sequence representing absent
576 genes, as well as a system for identifying misalignments by searching for erroneously-called
577 heterozygous SNP alleles in hemizygous genes. Thus, accurate copy number determination is a
578 vital first step in interpreting *KIR* sequencing data. To achieve this goal, we developed a copy

579 number determination method in PING that uses all described *KIR* alleles as an alignment resource,
580 increasing the total number of reference sequences from 15 to 905 compared to single-sequence
581 per gene alignments. While many of these alleles were only defined across exonic regions,
582 rendering them ineffective for short-read alignment, we developed and implemented a protocol for
583 intronic region imputation. The imputation method cannot resolve all uncharacterized nucleotide
584 sequence, yet it accounts for the majority of missing sequence, greatly increasing the number of
585 useful reference alleles. The exhaustive alignment provides a comprehensive map of the alleles to
586 which a read may align, facilitating copy number resolution of important *KIR* allelic groups and
587 genes that share extensive sequence similarity, such as *KIR2DL2*, *KIR2DL3*, *KIR2DL1* and
588 *KIR2DS1*, which were inaccessible to previous bioinformatic methods (27,44). Additionally, the
589 limited range of described UTR sequence, ~250bp 5'UTR and ~500bp 3'UTR, can reduce
590 alignments over the first exon and potential regulatory regions (45,46).

591

592 The improved copy determination performance of PING, in addition to the expanded useful
593 reference sequence repertoire, enables a smart, genotype-aware alignment workflow, designed to
594 minimize read misalignments by closely matching reference sequences to the gene sequences
595 present in the sequencing data. This alignment strategy addresses a major weakness of the filtration
596 alignments utilized in the prototype workflow, which apply filters to retain gene-specific reads and
597 eliminate cross-mapping reads regardless of the gene-content or sequence makeup, and thus often
598 suffer from either inadequate or patchy aligned read depths after filtration. The genotype-aware
599 workflow achieves highly accurate genotype determinations for the European dataset (Table 2),
600 and highly accurate resolved genotype determinations across all tested datasets (Table 3).

601

602 Both the synthetic dataset and Khoisan cohort showed higher levels of unresolved genotypes
603 compared to the European cohort (Table 2). These datasets represent challenging data to correctly
604 interpret, with the Khoisan being an extremely divergent population with many unresolved
605 genotypes in the validation data, and the synthetic dataset consisting of random alleles, some of
606 which used imputed sequence. An analysis into the discordant results for the synthetic dataset
607 (Figure 6, S8 Table) showed a complex web of cross-mapped reads. These cross-mapped reads
608 can be extremely difficult to resolve because the high-degree of sequence shared among *KIR* genes
609 (Figure 4) makes it almost impossible to determine correct mappings. Additionally, measures
610 meant to prevent read misalignments, such as the use of virtual probes to refine reference sequence
611 selection, can serve as a double-edged sword, where the issue at hand is addressed but the changes
612 create new sources for read misalignments. We believe improved interpretation of *KIR* sequencing
613 data will ultimately be achieved through longer-range sequencing technologies that can extend
614 past the range of the shared sequence motifs, and through better imputation approaches that can
615 more fully characterize currently described *KIR* alleles to provide a more robust alignment
616 reference. Meanwhile, for samples for which genotypes are not easily resolvable, we recommend
617 direct visualization of sequence alignments potentially coupled with alternative laboratory
618 methods to more precisely determine genotypes.

619

620 An analysis into the discordant copy results highlights a major outstanding problem with the PING
621 workflow since accurate copy determination is a central component of effective genotype-aware
622 alignments, and the need for manual thresholding between copy groups introduces the component
623 of user error. Continued development of the pipeline will address methods for automating copy

624 determination for targeted sequencing data that matches or surpasses the accuracy achieved by
625 manual thresholding.

626

627 While the PING workflow is specific to interpreting sequence originating from the *KIR* complex,
628 the underlying strategies can be extended to other problematic genomic regions. For example,
629 multiple-sequence per gene alignment strategies provide information for discriminating between
630 reads derived from genes with high sequence identity and extensive nucleotide polymorphisms.
631 Additionally, genotype-aware alignment strategies reduce bias introduced by the reference
632 sequence for reads derived from genomic regions with high structural variation.

633

634 In conclusion, PING incorporates these innovations to provide accurate, high-throughput
635 interpretation of the *KIR* region from short-read sequencing data. Together, these modifications
636 provide a consistent *KIR* genotyping pipeline, creating a highly automated, robust workflow for
637 interpreting *KIR* sequencing data. To the best of our knowledge, this is the only bioinformatic
638 workflow currently available for high-resolution *KIR* genotyping from short-read data. Given the
639 importance of *KIR* variation in human health and disease, availability of a highly accurate method
640 to assess *KIR* genotypic variation should promote important discoveries related to this complex
641 genomic region.

642

643 **Supporting information**

644 **S1 Figure. European cohort copy determinations.** (PDF)

645 **S2 Figure. Khoisan cohort copy determinations.** (PDF)

646 **S3 Figure. Synthetic dataset copy determinations.** (PDF)

- 647 **S1 Table. PING determined copy number table. (XLSX)**
- 648 **S2 Table. Validation copy number table. (XLSX)**
- 649 **S3 Table. PING determined genotype table. (XLSX)**
- 650 **S4 Table. Validation genotype table. (XLSX)**
- 651 **S5 Table. Virtual probe table for reference modifications. (XLSX)**
- 652 **S6 Table. K-mer gene match table. (XLSX)**
- 653 **S7 Table. Diverse and minimized reference allele set. (XLSX)**
- 654 **S8 Table. Discordant genotype analysis for synthetic dataset. (XLSX)**
- 655 **S1 Text. Genotype determination supporting methods. (DOCX)**

656

657 **Author contributions**

658 **Conceptualization:** Wesley Marin, Paul J. Norman, Jill A. Hollenbach

659 **Data Curation:** Wesley M. Marin, Ravi Dandekar, Danillo G. Augusto, Tasneem Yusufali,
660 Bianca Heyn, Jan Hoffmann, Vinzenz Lange, Jürgen Sauter, Paul J. Norman

661 **Formal analysis:** Wesley M. Marin

662 **Funding acquisition:** Jill A. Hollenbach

663 **Investigation:** Wesley M. Marin, Paul J. Norman

664 **Methodology:** Wesley M. Marin, Paul J. Norman

665 **Project Administration:** Jill A. Hollenbach

666 **Resources:** Paul J. Norman, Jill A. Hollenbach

667 **Software:** Wesley M. Marin, Ravi Dandekar, Tasneem Yusufali

668 **Supervision:** Jill A. Hollenbach

669 **Validation:** Wesley M. Marin, Paul J. Norman

670 **Visualization:** Wesley M. Marin

671 **Writing – Original Draft Preparation:** Wesley M. Marin

672 **Writing – Review & Editing:** Danillo G. Augusto, Paul J. Norman, Jill A. Hollenbach

673

674 **Acknowledgements**

675 We would like to thank Michael Wilson and Mark Seielstad for constructive comments. This work
676 was supported by the National Institutes of Health (NIH-R01AI128775).

677

678 **Funding**

679 This work was supported by the National Institutes of Health (NIH-R01AI128775). The funders
680 had no roles in study design, data collection and analysis, decision to publish, or preparation of the
681 manuscript.

682

683 **Conflict of Interest**

684 The authors have declared that no competing interests exist.

685

686 **References**

- 687 1. Colonna M, Moretta A, Vély F, Vivier E. A high-resolution view of NK-cell receptors:
688 Structure and function. In: Immunology Today. Elsevier Ltd; 2000. p. 428–31.
- 689 2. Björkström NK, Béziat V, Cichocki F, Liu LL, Levine J, Larsson S, et al. CD8 T cells
690 express randomly selected KIRs with distinct specificities compared with NK cells.
691 Blood. 2012 Oct 25;120(17):3455–65.

- 692 3. Fauriat C, Ivarsson MA, Ljunggren HG, Malmberg KJ, Michaëlsson J. Education of
693 human natural killer cells by activating killer cell immunoglobulin-like receptors. *Blood*.
694 2010 Feb 11;115(6):1166–74.
- 695 4. Pende D, Falco M, Vitale M, Cantoni C, Vitale C, Munari E, et al. Killer Ig-like receptors
696 (KIRs): Their role in NK cell modulation and developments leading to their clinical
697 exploitation. Vol. 10, *Frontiers in Immunology*. Frontiers Media S.A.; 2019. p. 1179.
- 698 5. Kumar S. Natural killer cell cytotoxicity and its regulation by inhibitory receptors
699 [Internet]. Vol. 154, *Immunology*. Blackwell Publishing Ltd; 2018 [cited 2020 Mar 5]. p.
700 383–93. Available from: <http://doi.wiley.com/10.1111/imm.12921>
- 701 6. Parham P. MHC class I molecules and KIRs in human history, health and survival. *Nat*
702 *Rev Immunol* [Internet]. 2005 Mar [cited 2015 Feb 10];5(3):201–14. Available from:
703 <http://www.ncbi.nlm.nih.gov/pubmed/15719024>
- 704 7. Nelson GW, Martin MP, Gladman D, Wade J, Trowsdale J, Carrington M. Cutting Edge:
705 Heterozygote Advantage in Autoimmune Disease: Hierarchy of Protection/Susceptibility
706 Conferred by HLA and Killer Ig-Like Receptor Combinations in Psoriatic Arthritis. *J*
707 *Immunol*. 2004 Oct 1;173(7):4273–6.
- 708 8. Salie M, Daya M, Möller M, Hoal EG. Activating KIRs alter susceptibility to pulmonary
709 tuberculosis in a South African population. *Tuberculosis*. 2015 Dec 1;95(6):817–21.
- 710 9. Hirayasu K, Ohashi J, Kashiwase K, Hananantachai H, Naka I, Ogawa A, et al.
711 Significant association of KIR2DL3-HLA-C1 combination with cerebral malaria and
712 implications for co-evolution of KIR and HLA. *PLoS Pathog*. 2012 Mar;8(3).
- 713 10. Khakoo SI, Thio CL, Martin MP, Brooks CR, Gao X, Astemborski J, et al. HLA and NK
714 cell inhibitory receptor genes in resolving hepatitis C virus infection. *Science* (80-). 2004

- 715 Aug 6;305(5685):872–4.
- 716 11. Martin MP, Gao X, Lee JH, Nelson GW, Detels R, Goedert JJ, et al. Epistatic interaction
717 between KIR3DS1 and HLA-B delays the progression to AIDS. *Nat Genet.*
718 2002;31(4):429–34.
- 719 12. Giebel S, Locatelli F, Lamparelli T, Velardi A, Davies S, Frumento G, et al. Survival
720 advantage with KIR ligand incompatibility in hematopoietic stem cell transplantation from
721 unrelated donors. *Blood [Internet].* 2003 Apr 3 [cited 2019 Oct 11];102(3):814–9.
722 Available from: <http://www.bloodjournal.org/cgi/doi/10.1182/blood-2003-01-0091>
- 723 13. Cooley S, Trachtenberg E, Bergemann TL, Saeteurn K, Klein J, Le CT, et al. Donors with
724 group B KIR haplotypes improve relapse-free survival after unrelated hematopoietic cell
725 transplantation for acute myelogenous leukemia. *Blood [Internet].* 2009 Jan 15 [cited 2019
726 Oct 11];113(3):726–32. Available from:
727 [https://ashpublications.org/blood/article/113/3/726/25172/Donors-with-group-B-KIR-](https://ashpublications.org/blood/article/113/3/726/25172/Donors-with-group-B-KIR-haplotypes-improve)
728 [haplotypes-improve](https://ashpublications.org/blood/article/113/3/726/25172/Donors-with-group-B-KIR-haplotypes-improve)
- 729 14. Venstrom JM, Pittari G, Gooley TA, Chewing JH, Spellman S, Haagenson M, et al.
730 HLA-C-dependent prevention of leukemia relapse by donor activating KIR2DS1. *N Engl*
731 *J Med [Internet].* 2012 Aug 30 [cited 2019 Oct 11];367(9):805–16. Available from:
732 <http://www.ncbi.nlm.nih.gov/pubmed/22931314>
- 733 15. Nakamura R, Gendzekhadze K, Palmer J, Tsai NC, Mokhtari S, Forman SJ, et al.
734 Influence of donor KIR genotypes on reduced relapse risk in acute myelogenous leukemia
735 after hematopoietic stem cell transplantation in patients with CMV reactivation. *Leuk Res.*
736 2019 Dec 1;87:106230.
- 737 16. Cooley S, Weisdorf DJ, Guethlein LA, Klein JP, Wang T, Marsh SGE, et al. Donor Killer

- 738 Cell Ig-like Receptor B Haplotypes, Recipient HLA-C1, and HLA-C Mismatch Enhance
739 the Clinical Benefit of Unrelated Transplantation for Acute Myelogenous Leukemia. *J*
740 *Immunol.* 2014 May 15;192(10):4592–600.
- 741 17. Guethlein LA, Beyzaie N, Nemat-Gorgani N, Wang T, Ramesh V, Marin WM, et al.
742 Following transplantation for AML, donor KIR Cen B02 better protects against relapse
743 than KIR Cen B01. *J Immunol.* Forthcoming.
- 744 18. Martin AM, Freitas EM, Witt CS, Christiansen FT. The genomic organization and
745 evolution of the natural killer immunoglobulin-like receptor (KIR) gene cluster.
746 *Immunogenetics* [Internet]. 2000 Apr 4 [cited 2019 Oct 11];51(4–5):268–80. Available
747 from: <http://www.ncbi.nlm.nih.gov/pubmed/10803839>
- 748 19. Uhrberg M, Valiante NM, Shum BP, Shilling HG, Lienert-Weidenbach K, Corliss B, et al.
749 Human diversity in killer cell inhibitory receptor genes. *Immunity.* 1997 Dec 1;7(6):753–
750 63.
- 751 20. Wilson MJ, Torkar M, Haude A, Milne S, Jones T, Sheer D, et al. Plasticity in the
752 organization and sequences of human KIR/ILT gene families. *Proc Natl Acad Sci*
753 [Internet]. 2000 Apr 25 [cited 2019 Oct 11];97(9):4778–83. Available from:
754 <http://www.ncbi.nlm.nih.gov/pubmed/10781084>
- 755 21. Robinson J, Halliwell JA, McWilliam H, Lopez R, Marsh SGE. IPD--the Immuno
756 Polymorphism Database. *Nucleic Acids Res* [Internet]. 2013 Jan [cited 2015 Feb
757 13];41(Database issue):D1234-40. Available from:
758 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3531162&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3531162&tool=pmcentrez&rendertype=abstract)
759 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3531162&tool=pmcentrez&rendertype=abstract)
- 760 22. Norman PJ, Abi-Rached L, Gendzekhadze K, Hammond JA, Moesta AK, Sharma D, et al.

- 761 Meiotic recombination generates rich diversity in NK cell receptor genes, alleles, and
762 haplotypes. *Genome Res.* 2009 May 1;19(5):757–69.
- 763 23. Traherne JA, Martin M, Ward R, Ohashi M, Pellett F, Gladman D, et al. Mechanisms of
764 copy number variation and hybrid gene formation in the KIR immune gene complex.
765 [cited 2020 May 6]; Available from: <http://www.ebi.ac.uk/ipd/kir/>
- 766 24. Hollenbach JA, Nosedal I, Ladner MB, Single RM, Trachtenberg EA. Killer cell
767 immunoglobulin-like receptor (KIR) gene content variation in the HGDP-CEPH
768 populations. *Immunogenetics* [Internet]. 2012 Oct 1 [cited 2019 Oct 11];64(10):719–37.
769 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/22752190>
- 770 25. Hollenbach JA, Augusto DG, Alaez C, Bubnova L, Fae I, Fischer G, et al. 16(th) IHIW:
771 population global distribution of killer immunoglobulin-like receptor (KIR) and ligands.
772 *Int J Immunogenet* [Internet]. 2013 Feb [cited 2019 Oct 11];40(1):39–45. Available from:
773 <http://www.ncbi.nlm.nih.gov/pubmed/23280119>
- 774 26. Hsu KC, Chida S, Geraghty DE, Dupont B. The killer cell immunoglobulin-like receptor
775 (KIR) genomic region: Gene-order, haplotypes and allelic polymorphism. Vol. 190,
776 *Immunological Reviews*. 2002. p. 40–52.
- 777 27. Norman PJ, Hollenbach JA, Nemat-Gorgani N, Marin WM, Norberg SJ, Ashouri E, et al.
778 Defining KIR and HLA Class I Genotypes at Highest Resolution via High-Throughput
779 Sequencing. *Am J Hum Genet* [Internet]. 2016 Aug 4 [cited 2017 Oct 10];99(2):375–91.
780 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/27486779>
- 781 28. Roe D, Vierra-Green C, Pyo C-W, Eng K, Hall R, Kuang R, et al. Revealing complete
782 complex KIR haplotypes phased by long-read sequencing technology. *Genes Immun*
783 [Internet]. 2017 [cited 2019 Oct 11];18(3):127–34. Available from:

- 784 <http://www.ncbi.nlm.nih.gov/pubmed/28569259>
- 785 29. Uhrberg M, Parham P, Wernet P. Definition of gene content for nine common group B
786 haplotypes of the Caucasoid population: KIR haplotypes contain between seven and
787 eleven KIR genes. *Immunogenetics*. 2002;54(4):221–9.
- 788 30. Amorim LM, Santos THS, Hollenbach JA, Norman PJ, Marin WM, Dandekar R, et al.
789 Cost-effective and fast KIR gene-content genotyping by multiplex melting curve analysis.
790 HLA [Internet]. 2018 Dec 1 [cited 2021 Mar 11];92(6):384–91. Available from:
791 <https://pubmed.ncbi.nlm.nih.gov/30468002/>
- 792 31. Wagner I, Schefzyk D, Pruschke J, Schöfl G, Schöne B, Gruber N, et al. Allele-Level KIR
793 Genotyping of More Than a Million Samples: Workflow, Algorithm, and Observations.
794 *Front Immunol* [Internet]. 2018 Dec 4 [cited 2019 Oct 11];9:2843. Available from:
795 <http://www.ncbi.nlm.nih.gov/pubmed/30564239>
- 796 32. Anderson KM, Augusto DG, Dandekar R, Shams H, Zhao C, Yusufali T, et al. Killer Cell
797 Immunoglobulin-like Receptor Variants Are Associated with Protection from Symptoms
798 Associated with More Severe Course in Parkinson Disease. *J Immunol* [Internet]. 2020
799 Sep 1 [cited 2021 Mar 11];205(5):1323–30. Available from:
800 <https://www.jimmunol.org/content/205/5/1323>
- 801 33. Vargas L de B, Dourado RM, Amorim LM, Ho B, Calonga-Solís V, Issler HC, et al.
802 Single Nucleotide Polymorphism in KIR2DL1 Is Associated With HLA-C Expression in
803 Global Populations. *Front Immunol* [Internet]. 2020 Aug 21 [cited 2021 Mar 11];11.
804 Available from: </pmc/articles/PMC7478174/>
- 805 34. Gu Z, Gu L, Eils R, Schlesner M, Brors B. circlize implements and enhances circular
806 visualization in R. *Bioinformatics* [Internet]. 2014 Oct 1 [cited 2019 Oct 11];30(19):2811–

- 807 2. Available from: <https://academic.oup.com/bioinformatics/article->
808 [lookup/doi/10.1093/bioinformatics/btu393](https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btu393)
- 809 35. R Core Team. R: A Language and Environment for Statistical Computing [Internet].
810 Vienna, Austria; 2018. Available from: <https://www.r-project.org/>
- 811 36. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
812 Alignment/Map format and SAMtools. *Bioinformatics* [Internet]. 2009 Aug [cited 2021
813 Mar 11];25(16):2078–9. Available from: </pmc/articles/PMC2723002/>
- 814 37. Leaton LA, Shortt J, Kichula KM, Tao S, Nemat-Gorgani N, Mentzer AJ, et al.
815 Conservation, extensive heterozygosity, and convergence of signaling potential all
816 indicate a critical role for KIR3DL3 in higher primates. *Front Immunol*. 2019;10(JAN).
- 817 38. Huang W, Li L, Myers JR, Marth GT. ART: A next-generation sequencing read simulator.
818 *Bioinformatics* [Internet]. 2012 Feb [cited 2021 Mar 11];28(4):593–4. Available from:
819 </pmc/articles/PMC3278762/>
- 820 39. Jiang W, Johnson C, Jayaraman J, Simecek N, Noble J, Moffatt MF, et al. Copy number
821 variation leads to considerable diversity for B but not A haplotypes of the human KIR
822 genes encoding NK cell receptors. *Genome Res*. 2012 Oct;22(10):1845–54.
- 823 40. Hollenbach JA, Norman PJ, Creary LE, Damotte V, Montero-Martin G, Caillier S, et al. A
824 specific amino acid motif of HLA-DRB1 mediates risk and interacts with smoking history
825 in Parkinson’s disease. *Proc Natl Acad Sci U S A* [Internet]. 2019 Apr 9 [cited 2021 Mar
826 11];116(15):7419–24. Available from: www.pnas.org/cgi/doi/10.1073/pnas.1821778116
- 827 41. Nemat-Gorgani N, Guethlein LA, Henn BM, Norberg SJ, Chiaroni J, Sikora M, et al.
828 Diversity of KIR, HLA Class I, and Their Interactions in Seven Populations of Sub-
829 Saharan Africans. *J Immunol*. 2019 May 1;202(9):2636–47.

- 830 42. Marin WM, Dandekar R, Augusto DG, Yusufali T, Norman PJ, Hollenbach JA. PING
831 [Internet]. Github. 2020 [cited 2020 Sep 29]. Available from:
832 <https://github.com/wesleymarin/PING>
- 833 43. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods
834 [Internet]. 2012 Apr 4 [cited 2019 Oct 11];9(4):357–9. Available from:
835 <http://www.ncbi.nlm.nih.gov/pubmed/22388286>
- 836 44. Pyke RM, Genolet R, Harari A, Coukos G, Gfeller D, Carter H. Computational KIR copy
837 number discovery reveals interaction between inhibitory receptor burden and survival. Pac
838 Symp Biocomput [Internet]. 2019 [cited 2019 Oct 11];24:148. Available from:
839 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6417817/>
- 840 45. Li H, Wright PW, McCullen M, Anderson SK. Characterization of KIR intermediate
841 promoters reveals four promoter types associated with distinct expression patterns of KIR
842 subtypes. Genes Immun [Internet]. 2016 Jan 1 [cited 2021 Mar 11];17(1):66–74.
843 Available from: </pmc/articles/PMC4724278/>
- 844 46. Nutalai R, Gaudieri S, Jumnainsong A, Leelayuwat C. Regulation of KIR3DL3 expression
845 via mirna. Genes (Basel) [Internet]. 2019 Aug 1 [cited 2021 Mar 11];10(8). Available
846 from: </pmc/articles/PMC6723774/>
847
848