

Check for updates



## kir-mapper: A Toolkit for Killer-Cell Immunoglobulin-Like Receptor (KIR) Genotyping From Short-Read **Second-Generation Sequencing Data**

Erick C. Castelli<sup>1,2</sup> D | Raphaela Neto Pereira<sup>2</sup> | Gabriela Sato Paes<sup>2</sup> | Heloisa S. Andrade<sup>3</sup> | Marcel Rodrigues Ferreira<sup>2</sup> | Ícaro Scalisse de Freitas Santos<sup>2</sup> | Nicolas Vince<sup>4</sup> | Nicholas R. Pollock<sup>5,6</sup> | Paul J. Norman<sup>5,6</sup> | Diogo Meyer<sup>3</sup>

<sup>1</sup>Department of Pathology, School of Medicine, São Paulo State University (Unesp), Botucatu, Brazil | <sup>2</sup>Molecular Genetics and Bioinformatics Laboratory (GeMBio) - Experimental Research Unit, School of Medicine, São Paulo State University (Unesp), Botucatu, Brazil | 3Department of Genetics and Evolutionary Biology, Institute of Biosciences, University of São Paulo, São Paulo, Brazil | 4Center for Research in Transplantation and Translational Immunology, Nantes Université, INSERM, Nantes, France | 5Department of Biomedical Informatics, University of Colorado School of Medicine, Aurora, Colorado, USA | 6Department of Immunology and Microbiology, University of Colorado School of Medicine, Aurora, Colorado, USA

Correspondence: Erick C. Castelli (erick.castelli@unesp.br)

Received: 21 November 2024 | Revised: 30 January 2025 | Accepted: 12 February 2025

Funding: This work was supported by Fundação de Amparo à Pesquisa do Estado de São Paulo, 2021/14851-9; Coordenação de Aperfeiçoamento de Pessoal de Nível Superior, 88881.879003/2023-01; CAPES-COFECUB, Me 1044/24; Conselho Nacional de Desenvolvimento Científico e Tecnológico, 307031/2022-5; NIH, R01AI128775.

#### **ABSTRACT**

Killer cell immunoglobulin-like receptors (KIRs) regulate natural killer (NK) cell responses by activating or inhibiting their functions. Genotyping KIR genes from short-read second-generation sequencing data remains challenging as cross-alignments among genes and alignment failure arise from gene similarities and extreme polymorphism. Several bioinformatics pipelines and programs, including PING and T1K, have been developed to analyse KIR diversity. We found discordant results among tools in a systematic comparison using the same dataset. Additionally, they do not provide SNPs in the context of the reference genome, making them unsuitable for whole-genome association studies. Here, we present kir-mapper, a toolkit to analyse KIR genes from short-read sequencing, focusing on detecting KIR alleles, copy number variation, as well as SNPs and InDels in the context of the hg38 reference genome. kir-mapper can be used with whole-genome sequencing (WGS), whole-exome sequencing (WES) and sequencing data generated after probe-based capture methods. It presents strategies for phasing SNPs and InDels within and among genes, reducing the number of ambiguities reported by other methods. We have applied kir-mapper and other tools to data from various sources (WGS, WES) in worldwide samples and compared the results. Using long-read data as a truth set, we found that WGS kir-mapper analyses provided more accurate genotype calls than PING and T1K. For WES, kir-mapper provides more accurate genotype calls than T1K for some genes, particularly highly polymorphic ones (KIR3DL3 and KIR3DL2). This comparison highlights that the choice of method has to be considered as a function of the available data type and the targeted genes. kir-mapper is available at the GitHub repository (https://github.com/erickcastelli/kir-mapper/).

### 1 | Introduction

Killer cell immunoglobulin-like receptors (KIRs) are a group of immunomodulatory receptors expressed on the cell surface of Natural Killer (NK) cells and subsets of T lymphocytes [1, 2].

These receptors modulate the activity of NK cell responses by activating or inhibiting cell effector activity. The ligand specificity and mode of function of each KIR is determined by its genetic sequence. NK cells can directly kill diseased cells, facilitating and speeding up defences against pathogens or secrete

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is

© 2025 The Author(s). HLA: Immune Response Genetics published by John Wiley & Sons Ltd.

HLA. 2025: 105:e70092

cytokines, such as interferon-gamma (IFN), which stimulate adaptive immune cell responses [3, 4]. KIR major ligands are the major histocompatibility complex (MHC) class I molecules, in humans termed HLA, expressed by somatic cells [5–7].

In humans, the KIR locus is located on chr19 at the leukocyte receptor complex (LRC), with 13 genes and two pseudogenes [7]. The genomic complexity of the LRC results from an evolutionary history involving duplications, intergenic recombination, point mutations and deletions resulting in a set of highly polymorphic genes; they display copy number variation and show high sequence similarities between paralogous copies [8–11]. The KIR genes are the most polymorphic receptors of human NK cells, displaying high allelic diversity [7, 11]. In addition to the presence of SNPs and InDels across all KIR genes, they also show copy number variation and can be described as absent or present from some KIR haplotypes, although the absence of *KIR3DL3* is a rare event [12]. The IPD-KIR database (version 2.13), an official repository for known KIR alleles, currently reports 2219 alleles (i.e., distinct sequences) considering all KIR genes [13].

Because KIR genes comprise a highly polymorphic multigene family, surveying KIR polymorphisms from short-read secondgeneration sequencing (NGS) data is challenging. Similarities between genes cause cross-alignments, with reads from one gene aligning to multiple loci [14]. In addition, high polymorphism may cause alignment failure when a single reference genome is used. Cross-alignment and alignment failure can lead to high error rates when genotyping KIR genes with conventional short-read aligners (e.g., BWA [15] and Bowtie2 [16]) and a single reference genome. Therefore, variant calls across KIR genes from most genome initiatives, such as the 1000Genomes project [17], may be biased or absent. HLA genes within the MHC at chromosome 6, which encode the major ligands for KIR receptors and present a similar organisation, face the same issue [18]. As for HLA, it is essential to use tools tailored to the LRC's structure and polymorphism to reliably genotype KIR genes using NGS short reads.

Some bioinformatics pipelines and programs have been developed to survey KIR allelic diversity. These include the original version of PING [11, 14], designed to genotype KIR from targeted sequencing using probes to specifically amplify KIR genes; T1K [19], compatible with RNA-seq and DNA-seq, including whole-exome sequencing (WES) and whole-genome sequencing (WGS); KIRCLE [20], which supports WES and uses blast to detect KIR alleles; and a recent PING update that supports WGS [21]. Currently, the most widely used KIR genotyping strategy is a well-established biotinylated DNA probe-based capture method coupled with PING [11, 14].

Here, we explore the performance of existing methods and compare these to a new toolkit to survey KIR genes, presented here for the first time. As expected for such a complex region, we observe conflicting results when we apply distinct bioinformatic tools to the same set of samples. An additional difficulty arises as available tools detect KIR alleles and sometimes copy numbers, but do not report SNPs and InDels in the context of the reference genome unless additional data post-treatment is applied. Therefore, their use for association studies is not straightforward.

We introduce kir-mapper, a toolkit for surveying KIR genes from short-read sequencing data. It focuses on detecting KIR alleles, copy number variation, SNPs and InDels across all KIR genes in the context of the hg38 reference genome. kir-mapper differs from existing KIR typing tools because it reports SNVs in the context of the hg38 reference genome and uses the inferred phase observed among these SNVs to define alleles. We have applied this method to survey KIR genes from different sources (WGS, WES and the probe-based capture method) in worldwide samples. We also compared kir-mapper, PING and T1K genotype calls. The rationale is that since these methods apply different algorithms for alignment, copy number determination and genotyping, the overlap and comparison between them would greatly support accurate KIR genotyping.

#### 2 | Methods

## 2.1 | The Kir-Mapper Workflow

kir-mapper is a toolkit designed to handle Illumina short reads. kir-mapper encompasses four primary functions embedded in a single program: *map*, for read alignment against the reference genome; *ncopy*, to detect copy number variation; *genotype*, for genotyping SNVs and InDels across all genes calling KIR alleles and *haplotype*, for phasing all variants, including those between genes, and resolving ambiguities not handled by genotype (Figure 1).

#### 2.2 | The Kir-Mapper Map Function

In the first step, a Kmer approach is used to identify reads that present at least 25 nucleotides matching any known KIR gene and sorts them into gene-specific fastq files. One read or pair of reads may be compatible with more than one gene. Then, using a motif approach to search for non-polymorphic sequences in each KIR gene, kir-mapper determines the presence or absence of each KIR gene. For the genes present in the sample, a scoring process calculates the distance (number of different nucleotides) between each read (or pair) and known KIR sequences. The known sequence database is composed of KIR alleles available on the IPD-KIR database [13], along with sequences from GENBANK [22] and sequences characterised in our lab from PCR amplification and Illumina sequencing. The distances are then compared to assign the reads to the most likely gene. Therefore, kir-mapper performs a multi-referenced alignment.

Despite their polymorphic nature, many KIR genes share high sequence similarities. Therefore, the software sometimes assigns a read to more than one gene. In these cases, kir-mapper treats all these alignments for these reads as secondary, which will be ignored in further steps. This occurs because some KIR genes share identical sequences in specific regions, thus explaining the presence of reads that do not provide unambiguous information about the locus they belong to. The map function also applies an algorithm that detects two sequences for each gene from the database of known KIR sequences that best fit the observed reads. Based on these detected sequences, the program recovers some secondary aligned reads, returning them to primary alignments when possible.

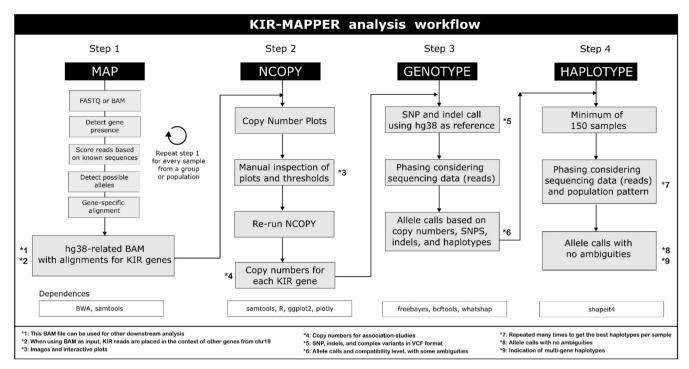


FIGURE 1 | The kir-mapper workflow to call KIR variants (SNVs and InDels), KIR alleles and haplotypes from short-read data.

After the scoring step, the program generates gene-specific fastq files, which are then aligned to the hg38 reference sequence from the gene associated with the fastq file. The alignments are combined, and the alignment positions are adjusted based on the hg38 reference genome. In summary, kir-mapper uses a multi-referenced alignment to reduce cross-alignments and alignment failures, which can occur with conventional short-read alignment tools and a single reference genome. The read alignment is subsequently adjusted to the coordinates of a single reference genome, producing a BAM file with reads aligned to the hg38 reference genome.

If a KIR gene is annotated in the primary assembly (hg38) of chr19, the reads mapping to this gene are aligned to chr19. However, if a gene is annotated in an alternative contig such as *KIR2DL5A* and *KIR2DL5B*, the reads will be aligned to only one of these alternative contigs. The list of genes annotated in alternative contigs and their positions is available at the GitHub repository (https://github.com/erickcastelli/kir-mapper). kir-mapper treats all *KIR* genes separately. This means that *KIR2DL2* and *KIR2DL3*, *KIR3DS1* and *KIR3DL1* and *KIR2DS1* and *KIR2DS4* are considered different genes with individual copy numbers and allele calls. The only exceptions are *KIR2DL5A* and *KIR2DL5B*, which are grouped as a single gene (KIR2DL5AB).

The map function requires raw fastq data (either paired or single-end) or a BAM file with reads aligned to the hg38 reference genome using BWA-MEM [15]. The output is a BAM file containing the reads aligned to the hg38 reference genome. This final BAM file can be examined using the Integrative Genome Viewer (IGV) and used in downstream analysis to identify SNVs across chr19 and alternative contigs.

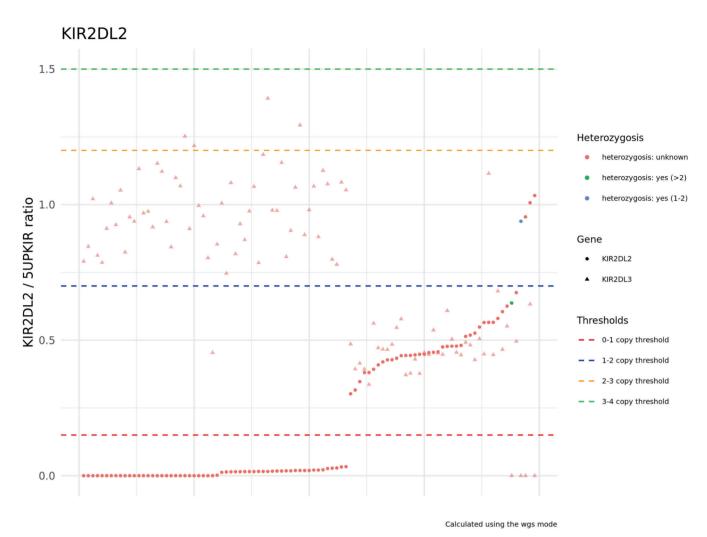
Under the hood, the map function utilises a database of known KIR alleles and relies on the BWA-MEM [15] algorithm to align

and calculate the distance between each read and the sequences from the database. Additionally, samtools is used to handle SAM and BAM files [23], and picard tools are used to mark duplicates. By default, the map function considers intronic sequences, which are suitable for WGS and the probe-based capture and sequencing. When processing WES, the proper flag (-exome) must be used to exclude intronic sequences from the map function.

#### 2.3 | The Kir-Mapper Ncopy Function

The kir-mapper ncopy is a tool that can detect the number of copies of all KIR genes. The input for ncopy is BAM files produced with the map function, and the user can process thousands of BAM files simultaneously. The tool uses samtools [23] to determine the depth of the gene regions and calculates the ratio between the depth of the target and the reference. *KIR3DL3* is the default reference since it is expected to be in two copies for most individuals [12], one per chr19. However, since all KIR genes may have copy number variation, as we will demonstrate here, the user may choose other alternative references such as *HLA-G*, *HLA-E* or a region upstream of the *KIR3DL3* gene between *ILT2* (*LILRB1*) and *KIR3DL3*, which we call 5UPKIR. This enables the user to test whether *KIR3DL3* is genuinely a framework gene for all samples.

The end products are gene-specific plots in .png and .html formats, indicating the ratio observed for each sample, as illustrated in Figure 2. The HTML version is an interactive plot. These plots were highly inspired by PING [11, 14]. By visually examining the plots, the user can identify the optimal thresholds, which are the points of transition between groups of samples with similar patterns. The user can modify these thresholds by editing a text file and running ncopy to recalculate copy numbers of all genes and samples (Figure 2).



**FIGURE 2** | A plot generated by kir-mapper ncopy, reporting the ratios between *KIR2DL2* depth and depth of a region between *KIR3DL3* and *ILT2* (here called 5UPKIR) for all Finnish samples from the 1000Genomes dataset. Individuals are ordered from left to right based on the estimated ratio. Samples can be assigned to different groups. In this case, there are three groups: Those with no *KIR2DL2* (ratios around zero); samples with one copy of *KIR2DL2* (ratios around 0.5); and samples with two copies (ratios around 1.0). The translucid triangles represent the ratios observed for *KIR2DL3*. Because *KIR2DL3* occupies the same genomic location as *KIR2DL2*, the presence of one copy is expected to preclude the presence of the other, and their occurrences should be complementary.

Determining the copy number is crucial for accurate genotyping. It is highly recommended that *KIR3DL3* be evaluated as a suitable reference. Upon examining the Finnish samples from the 1000 Genomes dataset [17], we identified two individuals, HG00273 and HG00378, who appear to have three copies of *KIR3DL3* when using the 5UPKIR, *HLA-G* or *HLA-E* as references. Therefore, *KIR3DL3* is not a suitable reference in Finland. We also detected one sample from the SABE/Brazil cohort [24] with 3 copies of *KIR3DL3*. Figure S1 is an example of the plot produced by ncopy for *KIR3DL3* using 5UPKIR as a reference, with all Finnish samples from the 1000Genomes project [17].

#### 2.4 | The Kir-Mapper Genotype Function

Genotyping of SNVs within all KIR genes is performed using freebayes [25] and relies on the results for copy numbers and alignments produced using map and ncopy functions from previous steps. kir-mapper contains an algorithm that removes

artefacts and classifies the uncertain variants with low depth or unbalanced heterozygotes as missing alleles. Next, whatshap [26] phases sites that are heterozygous in the sample and occur on the same read. This process is crucial in reducing the possible allele combinations for each gene, which ultimately minimises ambiguities. Finally, all the variants are reported in the context of the hg38 reference genome, resulting in partially phased gene-specific VCF files.

The comparison between the observed variants (SNVs and InDels) and the phasing of these variants within each KIR locus plus the patterns observed in known KIR alleles determines the most likely allele combination for each sample. This comparison takes into account both the variants observed in known KIR alleles and any new variants that might have been detected. The outcome of this process is a text report for each sample that includes the total number of tested variants, the proportion of matches and mismatches between the sample and the tested alleles and a list of any potential mismatches observed for the tested alleles.

Despite the efforts to detect the haplotypes within each KIR gene and resolve potential ambiguities, the genotype function may report a list of possible allele combinations that fit all the observed SNVs and microhaplotypes. This can happen because, when using short reads with whatshap [26], only variants within the same read or pair of reads can be phased. Therefore, not all heterozygous sites are phased, particularly when analysing WES data, in which variants in different exons are too far apart. The following section presents the haplotype function, which was designed to solve such ambiguities.

# 2.5 | The Kir-Mapper Haplotype Function—Solving Ambiguities

Ambiguities are common with KIR gene genotyping because several allele combinations might present identical SNPs and InDels. These SNPs must be phased into two or more haplotypes per gene to solve these ambiguities. However, phasing distant variants when dealing with short reads is sometimes impossible for approaches aiming to extract the phasing status directly from the sequencing reads. Ambiguities would not be an issue when using long reads, but kir-mapper is still incompatible with long reads. Although the list of possible allele combinations is highly reduced by the whatshap phasing step, which can theoretically phase all variants within an exon, ambiguities still occur when the program does not detect the phase between variants in adjacent exons.

kir-mapper has a built-in method to resolve ambiguities, but it requires the genotyping of at least 150 samples simultaneously. We use statistical phasing to assemble within a single haplotype the smaller haplotypes defined by whatshap [26]. Then, we define two haplotypes per gene. The first step is to convert the gene-specific VCFs into a dummy diploid VCF that includes dummy positions for all variants and genes, respecting the known order of the KIR genes. This diploid VCF also includes information on the presence or absence of each gene. All variant genotypes are re-encoded to fit the diploid state. This is done by adding a reference allele when there is a deletion for the gene or reducing triploid and tetraploid variants to diploid when there are only one or two possible alleles at the site. If it is not possible to reduce the genotype to only two alleles, the genotype is replaced by missing alleles. Then, we use shapeit4 [27] to phase all the variants. We run shapeit4 multiple times, and haplotypes are compared to select the haplotype that appears most often. This step preserves the phasing sets detected by whatshap when applying the genotype function. The outcome is a fully phased VCF.

Then, kir-mapper creates two sequences for each sample and each gene by using the reference genome sequence, all the SNVs detected and the phase among these SNVs. It then compares those sequences with the ones available in the IPD-KIR database [13]. The outcomes are two alleles per gene per sample, with no ambiguities. The absence of the gene is reported as allele \*null. The comparison between the calls from the genotype and haplotype functions assists in resolving ambiguities. For the haplotype function, users are cautioned with a message in the final report when haplotypes indicate the presence of three or four copies of a KIR gene because the

program will report only two alleles that might not reflect the true genotypes. In these cases, users should consider only the kir-mapper genotype calls.

# 2.6 | Testing Kir-Mapper With Sequencing Data From Multiple Sources

We tested kir-mapper in five different ways. First, we simulated Illumina HIseq 2500 sequencing data for 25 samples, with a 2  $\times$  150bp reads, a fragment size of 450±150 and a target depth of 60X, by using art\_illumina [28]. Each virtual sample presented two copies for each KIR gene, with a known allele for each of these copies. We acknowledge this is unrealistic, but this configuration creates the most difficult scenario for KIR genes alignment. We aligned reads using BWA-MEM [15] and hg38 as reference (with alternative contigs and HLA alleles; the same used by the 1000 Genomes project [17]) and with the kirmapper map function. We then tracked the gene where the read originated and where it aligned. We produced plots tracking the alignments by using R and ggplot2.

Second, we applied kir-mapper to 172 samples with KIR alleles called using the latest version of PING and the probe capture and sequencing method [11], from a study addressing KIR alleles and susceptibility to COVID-19 [29]. PING and this sequencing method are widely used in multiple studies of KIR diversity.

Third, we tested kir-mapper in 34 samples with long-read phased assemblies from the Human Pangenome Reference Consortium (HPRC) [30, 31], in which Illumina short-read data is also available [17], comparing the KIR allele calls reported for the long reads and obtained with kir-mapper and short reads.

Fourth, we applied kir-mapper to survey the KIR alleles in samples from the 1000 Genomes project [17], all with Illumina short reads and depth around 30X. For this, we selected one population with major ancestry from each biogeographic region: YRI (Yoruba in Ibadan, Nigeria), GBR (British from England and Scotland), CLM (Colombian in Medellín, Colombia), JPT (Japanese in Tokyo, Japan) and ITU (Indian Telugu in the UK). We genotyped KIR genes by using PING, T1K and kir-mapper, comparing the results. We downloaded the BAM files with reads aligned to the hg38 reference genome, which were used as input for kir-mapper. We also converted this data to fastq to be used with T1K and PING. For kir-mapper, we ran the map step for each sample from a specific population, using the name of the population as the output folder. This creates a kir-mapper output structure with all samples within the same population. Then, we ran the ncopy function for each population separately to determine copy numbers. Afterwards, we combined all samples (and populations) in a single folder using the kir-mapper function called "group" Finally, we ran the genotype and haplotype step, considering all samples simultaneously.

Fifth, we applied kir-mapper in the exome mode to evaluate the WES data for samples from the 1000 Genomes Project that presented the same genotype by PING, T1K and kir-mapper when evaluating the WGS data. The rationale of this analysis is to generate a truth set with well-documented KIR alleles and evaluate kir-mapper and T1K performances when evaluating WES data.

### 3 | Results and Discussion

# 3.1 | Alignment and SNP Genotyping Performance Using Simulated Sequencing Data

We simulated NGS short reads for 25 samples (see methods). For each pair of reads, we tracked which KIR gene originated the read and where it aligned under two protocols: using BWA-MEM and a single reference genome and with kir-mapper.

When using BWA-MEM, there is a high degree of cross-alignment, with reads from one locus aligning to another. These cross-alignments occur because of the sequence similarities among KIR genes and also because some of the KIR genes are not present on the chr19 reference genome. Consequently, reads from genes not represented at the main chr19 sequence still align with the main chr19 sequence (Figure 3, left panel). For instance, sequences from the KIR2DS5 gene, which is not present in the main chr19 sequence from hg38, align with KIR2DL1, KIR2DL3 and KIR2DS4. Likewise, sequences from KIR3DS1 align with KIR3DL1. Sequence similarity also leads to crossalignments between genes that are present in the main chr19 sequence from hg38, such as KIR2DL1 and KIR2DL3. Because of this, read depth is much higher than the number of simulated reads in some regions and lower in others.

For *KIR2DL4* and *KIR3DL3*, we observe a different scenario. Many reads are aligned elsewhere or not aligned at all, reducing depth throughout the genes. Some samples presented a read depth as low as 40% of that expected by the simulation in specific regions. In addition, there are cross-alignments in some regions, particularly between *KIR3DL3* and *KIR3DS1*. This pattern of cross-alignment is different from that observed elsewhere, with reads from *KIR3DL3* mostly cross-aligning with *KIR3DP1* [14].

Using kir-mapper for these same simulated reads, read depth is homogeneous across the locus and close to the expected value for most genes (red line), and cross-alignments are rare (Figure 3, right panel). However, some reads still align to more than one gene because the alleles in the sample present the same sequence in some regions for two different genes. Therefore, all the alignments regarding these reads are marked as secondary and disregarded by the genotyping algorithm, reducing depth in some regions.

The map function from kir-mapper significantly improves alignment accuracy in all KIR genes, although some misalignments still occur, particularly for *KIR2DL1*. Figure S2 illustrates the alignment pattern for KIR genes not present at the main chr19 sequence from the hg38 reference genome. These optimised alignments from the function map significantly impact copy number determination for all KIR genes and the accuracy of detecting specific InDels and SNPs across each gene (Figure 4). For instance, copy number determination based on depth would be significantly impaired if depth were calculated based on

the BWA –M EM alignments. In such a case, genes with cross-alignments (*KIR2DL3*, for instance), genotyping would be biased with many false-positive and false-negative variants.

We also tested genotyping accuracy using freebayes. To establish the ground truth, we used simulations to force the alignment of simulated short reads from a specific gene to the reference of that specific gene. Therefore, there were no alignment errors such as misalignments or cross-alignments. After that, we genotyped the simulated samples with freebayes to obtain the expected genotype, thus defining the ground truth for SNPs in an errorfree environment when there was no alignment error. Then, we used freebayes to genotype SNPs and InDels after aligning the same reads using two different methods: BWA-MEM and the reference genome, and with kir-mapper, comparing the results with the previous ground truth. Therefore, the only modification is the alignment method, while the genotyping strategy (freebayes) was the same (Figure 4). While genotyping is extremely biased when using BWA-MEM and the hg38 reference genome (Figure 4, top panel), after using kir-mapper, the majority of the genotypes are identical to the truth set (Figure 4, bottom panel), with errors mostly in intronic regions from *KIR2DL1* (Figure 3). The genotyping errors observed when using BWA-MEM are mostly related to misaligned reads leading to the detection of false heterozygous sites (Figure 4, red positions). Therefore, this simulation indicates that genotyping data from chr19 within KIR genes should be considered with caution unless some KIRspecific method was applied to detect such genotypes.

## 3.2 | Calling KIR Alleles When Capturing KIR With Probes

We applied kir-mapper to 172 samples with KIR alleles genotyped by the latest version of PING and the probe-capture and sequencing method to enrich KIR [11] from a study addressing KIR polymorphism and susceptibility to COVID-19 [29]. Ambiguities reported by PING were solved with the PHASE program [32]. We compared the final calls (with no ambiguities) between methods. The only exception was *KIR3DP1*, for which we compared the raw calls (with ambiguities) between the two methods.

There is an important overlap in results obtained by kir-mapper and PING, with both methods detecting the same alleles and copy numbers (Figure 5). Despite the high overall overlap of results between methods, depending on the gene, between 2% and 10% of the samples presented different calls (light grey, light blue and black). In addition, both methods failed to genotype some samples and genes (shades of grey). The differences between methods mostly relate to differences in calls for copy numbers between methods or differing phasing of SNPs when making allele calls. The relatively low sample size (N=172) might explain different results when solving ambiguities using probabilistic models such as PHASE and Shapeit4, particularly for the most polymorphic genes, KIR3DL2 and KIR3DL3. Evaluating which method is correct is only possible if we apply other techniques, such as long-read sequencing. Therefore, it is essential to use multiple methods to evaluate KIR copy numbers and alleles and manually check possible inconsistencies. This manual check might be a visual inspection of the BAM file (reads aligned to

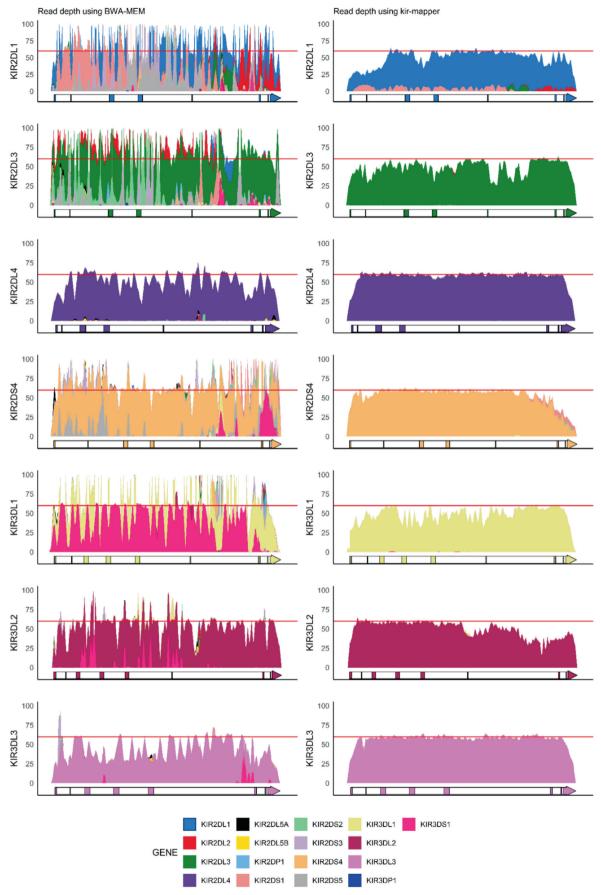


FIGURE 3 | The pattern of read alignment across seven KIR genes when using BWA-MEM and the reference genome hg38 (left panel) and when using kir-mapper (right panel). Different colours represent different origins for the reads. The gene structure is indicated below the x-axis, with boxes representing the exons. The horizontal red line represents the expected read depth (60X).

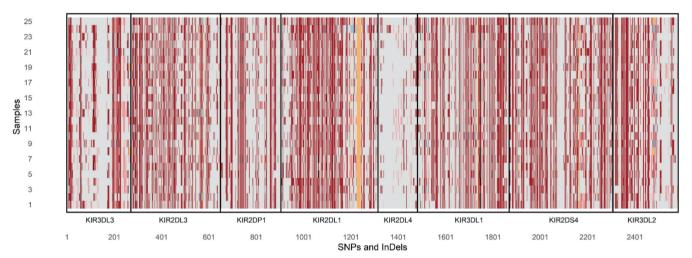
the hg38 reference genome) using IGV, for instance, which is only possible by using kir-mapper.

# 3.3 | Calling KIR Alleles From Whole-Genome Sequencing Data

Obtaining KIR genotypes from whole-genome sequencing offers a valuable source of data for studies of KIR diversity, evolution and genome structure across populations. Many publicly available whole-genome sequencing (WGS), including the 1000 Genomes Project [17], HGDP and SABE [24], can potentially be used to obtain well-curated and reliable KIR data. WGS also has advantages over WES data by avoiding probe bias and characterising intronic, regulatory and intergenic regions. However, as we will demonstrate, surveying KIR data from WGS is not easy, and available tools report different results for many samples.

We compared the allele calls from kir-mapper, PING-WGS and T1K with the ones reported for 34 long-read phased assemblies from the Human Pangenome Reference Consortium (HPRC) [30, 31], for which Illumina short-read data is also available from the 1000 Genomes project [17] (Figure 6). In our analyses, we treated the HPRC long-read calls as a truth set, and we compared the results of the analyses based on short-read data. kir-mapper performed better for all genes, followed by PING. The calls in which kir-mapper reported the same alleles as the long reads are marked in shades of blue, with accuracy varying from 89.2% for KIR3DP1 to 100% for KIR2DL2, KIR2DL4, KIR2DP1, KIR2DS1, KIR2DS3, KIR2DS5, KIR3DL2, KIR3DL3 and KIR3DS1. PING accuracy varied from 51.3% for KIR2DP1 to 100% for KIR2DS1, KIR2DS2 and KIR2DS5. T1K accuracy varied from 10% for KIR3DL3 to 97.3% for KIR2DS1. kir-mapper accuracy overcame PING and T1K for the most polymorphic genes, KIR3DL2 and KIR3DL3.

#### Genotyping accuracy using BWA-MEM + hg38



#### Genotyping accuracy using kir-mapper

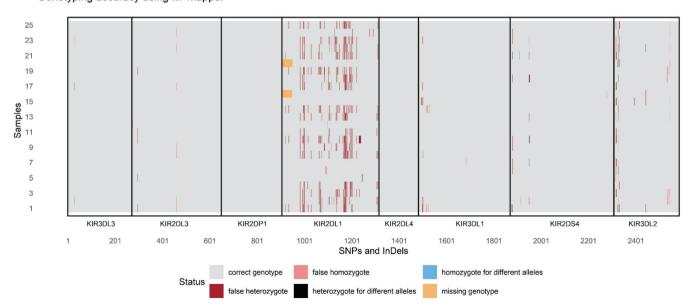


FIGURE 4 | KIR SNP and InDel genotyping accuracy when using BWA-MEM and the hg38 reference genome (upper panel) or when using the map function from kir-mapper to realign the reads to KIR genes. Genotyping was performed by freebayes, and the genotypes were compared to a truth set in simulated data. This simulation includes 25 samples with two random alleles for each KIR gene, a read size of 150 nucleotides, paired-end, read depth of 60× and 2584 SNPs or InDels.

## Kir-mapper vs PING comparison for KIR genotyping using probe-enriched Illumina sequencing

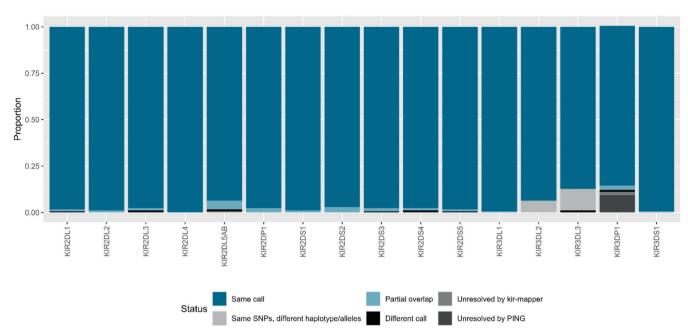


FIGURE 5 | Compatibility between kir-mapper and PING allele calls for KIR genes, for 172 samples sequenced using Illumina and the probecapture and enrichment technique. The comparison involved the categories: 'Same call' when the reported alleles are the same by both methods, 'Partial overlap' when one allele is identical and the other is different, 'Same SNPs, different haplotype/alleles' when both methods detected the same SNPs but the phasing process determined different alleles, 'Different call' when none of the alleles is the same by both methods, 'Unresolved by kir-mapper' when kir-mapper failed to report an allele combination, and 'Unresolved by PING' when PING failed to report an allele combination.

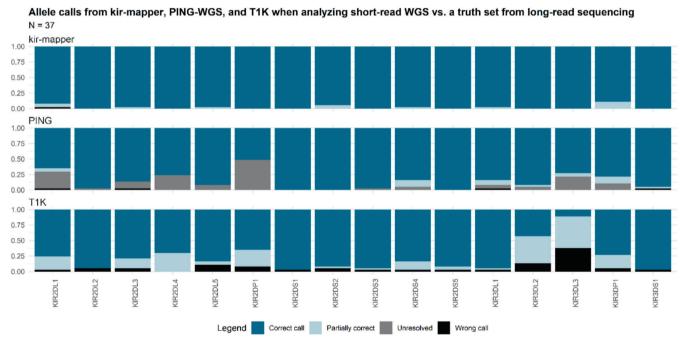


FIGURE 6 | Compatibility between KIR allele calls from kir-mapper, PING, and T1K using short reads and the alleles reported by phased assemblies from the Human Pangenome Reference Consortium (HPRC).

Next, we compared the final calls from kir-mapper (no ambiguities) with PING-WGS [14, 21] and from T1K [19] (Figure 7) for five populations from the 1000 Genomes, sampled from different continents. There was no truth set in this case, and our analyses focused on the degree of overlap across methods. Allele

calls from each method are available in Table S1. We found that the proportion of samples with the same call by all methods (dark blue) varies depending on the gene, ranging from 28.4% for *KIR3DL3* to 94% for *KIR2DS1*. There is also a high proportion of samples with different calls by all methods (median of

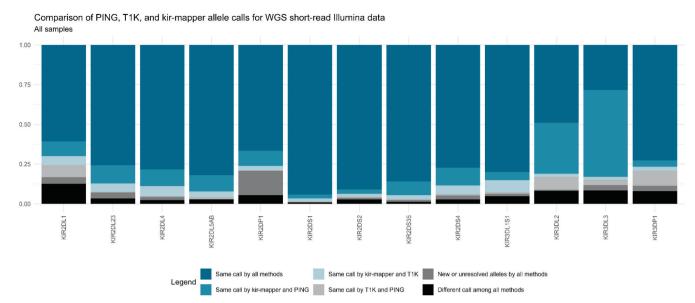


FIGURE 7 | Compatibility between kir-mapper, T1K and PING (the whole-genome version) allele calls for *KIR* genes when processing 30X whole-genome sequencing data for five populations from the 1000 genomes dataset. The comparison involved the categories: 'Same call by all methods' when the reported alleles are the same by all methods, 'Different call among all methods' when none of the alleles is the same by all the methods, 'New or unresolved alleles by all methods' when all methods agree that there is a new allele, 'Same call by T1K and kir-mapper' when the T1K and kir-mapper calls are compatible and PING reported a different one and 'Same call by kir-mapper and T1K', 'Same call by PING and T1K', 'Same call by kir-mapper and PING'.

6%, black). For KIR3DL3 and KIR3DL2, two of the most polymorphic KIR genes, PING and kir-mapper agree for most samples. However, T1K reported a different allele combination for almost half the samples. There is a higher compatibility between the calls from kir-mapper and PING and for most genes, except for KIR3DP1. We observed similar patterns when each population was evaluated separately but with a high proportion of samples with a different call by T1K as compared to kir-mapper and PING among African samples (Figure S3). The high proportion of samples with different genotypes by different methods highlights the difficulty of getting reliable genotypes for KIR genes. It became clear that each method's performance varies according to the KIR gene and the type of data being processed. While the high proportion of differences between methods is a cause for concern, the possibility of comparing short-read inferences to long-read results provides a criterion for establishing accuracy. As such, it indicates that the kir-mapper outperforms the other tested methods.

We monitored the frequencies of the alleles reported for all KIR genes and populations. In this case, we considered the kirmapper genotypes, even when this genotype is different from the ones reported by PING and T1K. We also removed samples with three or more gene copies to plot the allele frequencies (as shown in Figure 8 and Supporting Information). If a gene was absent on one chromosome, it was represented as allele \*null. Thus, an individual lacking *KIR2DL1*, for instance, has two copies of KIR2DL1\*null.

Figure 8 demonstrates the frequencies observed for all KIR2DL1 alleles detected among the five population samples. The occurrence of haplotypes without KIR2DL1 (the KIR2DL1\*null allele) is common in all populations but more prevalent among Europeans and populations with a significant European

ancestry, such as Colombians. The frequencies observed for *KIR2DL1* are compatible with those reported for other samples from the same biogeographic regions (www.allele-frequencies. net) [33]. For instance, the most prevalent *KIR2DL1* allele in East Asia is \*00302 (around 74%), as observed in the JPT group from the 1000 Genomes dataset. This same allele has a frequency of 36% in Ghana [34] and 40% in the YRI group, both from West Africa.

kir-mapper may report unresolved alleles when there are missing SNPs and Indels, leading to a long list of possible allele combinations. Unresolved alleles might occur due to low read depth or misalignments. kir-mapper also reports possible new alleles when none of the known alleles [13] match the observed genotypes. The proportion of new and unresolved alleles varies among KIR genes. The high proportion of new alleles in some KIR genes and populations, particularly in Africa, might reflect the underrepresentation of alleles from these populations in the IPD-KIR database [13] since their presence in the database depends on an accurate characterisation with a combination of long and shortread sequencing by NGS. We provide the allele frequencies of all KIR genes in Figure S4. The similarities observed among the frequencies reported for other populations from the same biogeographic region and those detected here are an encouraging indication that NGS analysis using kir-mapper provides results consistent with well-tested approaches of KIR typing.

## 3.4 | Calling KIR Alleles From Exomes

Exome data brings an additional challenge to KIR analysis. In addition to the cross-alignments and alignment failures, there is also probe bias, with one chromosome less captured than the other or not captured at all. kir-mapper was designed to

#### KIR2DL1

	KIR2DL1*00101	4.8	4.0	2.5	0.0	0.0	l
	KIR2DL1*00201	10.2	27.3	14.7	11.8	1.0	
Alleles	KIR2DL1*00302	47.3	33.0	40.2	73.5	39.5	
	KIR2DL1*00303	0.0	0.0	0.0	0.0	4.8	
	KIR2DL1*00306	0.0	0.0	0.0	0.0	0.5	
	KIR2DL1*00401	7.0	2.3	10.3	0.0	6.7	
	KIR2DL1*00601	0.0	0.0	0.0	0.0	4.8	
	KIR2DL1*007	0.5	0.0	0.0	0.0	4.3	Frequency (%)
	KIR2DL1*008	0.0	0.6	0.0	0.0	0.0	По
	KIR2DL1*010	1.6	2.3	2.5	0.0	1.4	10
	KIR2DL1*01102	0.0	0.0	0.0	0.0	3.3	
	KIR2DL1*01201	1.1	0.0	0.0	0.5	2.9	20
	KIR2DL1*01202	0.5	0.0	0.0	0.0	4.8	30
	KIR2DL1*014	0.0	0.0	0.0	0.0	2.4	40
	KIR2DL1*020	0.0	1.1	0.0	0.0	0.0	
₹	KIR2DL1*025	0.0	0.0	0.0	0.0	0.5	50
	KIR2DL1*03201	0.0	1.1	0.0	0.0	0.0	60
	KIR2DL1*034	0.5	0.6	0.5	1.5	0.0	70
	KIR2DL1*035	0.5	0.6	0.5	0.0	0.0	80
	KIR2DL1*03701	0.0	0.6	0.0	0.0	0.0	
	KIR2DL1*05101	3.2	4.0	2.5	0.0	1.0	90
	KIR2DL1*05401	0.5	1.1	0.0	0.0	1.0	100
	KIR2DL1*063	0.0	0.0	0.5	0.0	0.0	
	KIR2DL1*069	0.0	0.0	0.0	1.5	0.0	
	KIR2DL1*070	0.5	0.0	0.0	0.0	0.0	
	KIR2DL1*073	0.0	0.0	0.0	0.5	0.0	
	KIR2DL1*new	1.6	2.3	7.8	5.9	7.6	
	KIR2DL1*null	19.4	19.3	17.6	4.9	8.6	
ŀ	(IR2DL1*unresolved	0.5	0.0	0.5	0.0	5.2	
		CLM	GBR	ITU <b>Population</b>	JPT	YRI	

**FIGURE 8** | KIR2DL1 allele frequencies in five populations from the 1000 genomes dataset. KIR2DL1 was genotyped using kir-mapper. YRI (Yoruba in Ibadan, Nigeria, N=105), GBR (British from England and Scotland, N=88), CLM (Colombian in Medellín, Colombia, N=102), JPT (Japanese in Tokyo, Japan, N=102) and ITU (Indian Telugu in the UK, N=102). \*new alleles represent possible new alleles that are not in the IPD-KIR database.

determine copy numbers and genotypes from WGS and WES, using slightly different algorithms in each case.

Because we do not have samples with WES data and KIR genes validated by other methods to evaluate the kir-mapper (and T1K [19]) performance for WES, we opted for a different strategy to create a truth seq. We downloaded from the 1000 Genomes dataset the exome data from samples in which PING, T1K and kir-mapper called the same alleles when evaluating WGS (Figure 7, in dark blue) to be used as a truth set. Then, we applied T1K and kir-mapper to call KIR alleles, focusing only on the genes with 100% concordance between methods when evaluating the WGS data. Therefore, the sample size is different for each KIR gene. Afterwards, we compared the outputs obtained from the WES with those from the WGS (Figure 9).

The comparison demonstrated that the compatibility between the outputs when evaluating WGS or WES depends on the gene and the method. For most genes, the majority of the samples gave the same alleles using each method when evaluating WES, and these alleles were compatible with the

WGS data. The only exception was KIR2DL5A/B, in which T1K failed to detect the correct allele combination from the WES data in most samples. In addition, while both methods are equally efficient for some genes (KIR2DS4), kir-mapper or T1K performed better for others. T1K better detected alleles from WES data for KIR2DL1, KIR2DL2, KIR2DS1, KIR2DS2, KIR2DS5 and KIR3DL1. kir-mapper had a better performance for KIR2DL3, KIR2DL4, KIR2DL5A/B, KIR2DP1, KIR2DS3, KIR3DL2, KIR3DL3, KIR3DP1 and KIR3DS1. Regarding the two most polymorphic KIR genes, kir-mapper is more accurate for KIR3DL2, particularly among Europeans, Colombians, Japanese and Indian Telugu (Figure S5). For KIR3DL3, kirmapper performed better for Colombians and Indian Telugu, and T1K was better among Africans and Europeans. However, these results might be biased since the number of samples included for KIR2DL1, KIR2DL3, KIR2DL4, KIR2DP1, KIR2DS1, KIR3DL2, KIR3DL3 and KIR3DP1 is low because only a few samples had the same call by T1K, PING and kir-mapper in the WGS data (Figure 7). It should be noted that kir-mapper outperformed T1K for all genes when calling alleles in WGS data (Figure 6).

Compatibility among the T1K and kir-mapper calls from Exomes as compared to the WGS data.

The truth set includes all samples with the same call from T1K, kir-mapper, and PING from WGS data.

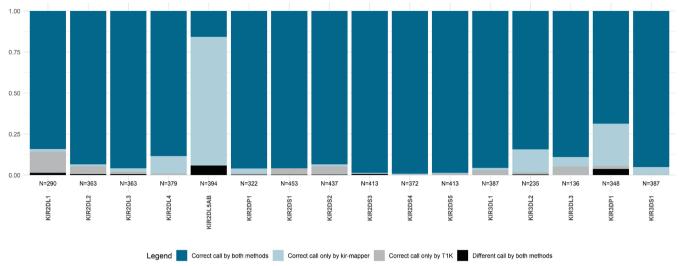


FIGURE 9 | Compatibility between the KIR genotypes detected by T1K and kir-mapper in whole-exome (WES) data with the ones reported when analysing whole-genome (WGS) data. WGS results were validated by three methods, T1K, PING-WGS and kir-mapper (Figure 7). Our comparison involved the categories: 'same call by both methods' when the alleles reported at the WES and WGS data are the same for both methods, 'different call by both methods' when the reported alleles are different from the WGS data and 'same call only by kir-mapper' when kir-mapper reported the same alleles when analysing WGS or WES data, but not T1K, and 'Same call only by T1K' when T1K reported the same alleles when analysing WGS or WES data, but not kir-mapper.

## 3.5 | Pros and Cons of PING, T1K and Kir-Mapper

T1K [19] is an easy-to-use and install tool, with some necessary customisation regarding the resolution level to be reported. It is a high-speed tool and demands minimal resources, and it is compatible with Linux and macOS. T1K processes individual samples, and it is not easy to evaluate different ploidies. For instance, if the output indicates just one allele for a KIR gene, users must decide based on the observed depth if they have one or more than one copy of that allele. Likewise, if the output indicates two different alleles for a KIR gene, users must decide based on the observed depths if both alleles are truly present (one might be an error, with very low depths) and the number of copies of each allele. Therefore, getting ploidy using T1K is not straightforward. We also noticed a poor performance of T1K for KIR3DL2, KIR3DL3 and KIR2DL5A/B when analysing WGS or WES data (Figures 6 and 9). Although T1K reports new SNPs on a VCF file, these reports are not in the context of the hg38 reference genome, and known SNPs are not reported automatically.

The PING version for WGS [14, 21] is slower than T1K and kirmapper and demands higher memory and processing power resources. It runs under R and RStudio and is dependent on specific program versions (likewise kir-mapper). It is designed to be run using a Linux server, and we were unable to install and run it on macOS. Once installed, it runs smoothly with samples from one population (about 100) in a 64Gb RAM machine, but it has crashed when running all the samples included in this study simultaneously due to lack of memory. PING outputs are easy to interpret. Unlike T1K, PING allows the analyses of several samples simultaneously and uses all the samples to determine copy numbers, influencing the genotyping process. Therefore, determining copy numbers with PING is an easy task. PING

reports all the variants observed across each KIR gene, but postprocessing is needed to place the SNPs in a VCF-like format, and these variants are not in the context of the hg38 reference genome. Compared to kir-mapper and T1K, PING reports more unresolved alleles. PING currently has no built-in tool to solve ambiguities for final allele calls, and most studies rely on inferring haplotypes with PHASE [32]. PING reports the allele combinations based on the copy numbers. Therefore, different from T1K, PING does not demand the manual evaluation of whether an allele with low read depth is in the sample. However, the current version of PING demands a large number of samples to accurately determine the copy numbers and genotypes of the samples, while T1K runs a single sample. PING and kirmapper results overlap better than T1K and PING or T1K and kir-mapper. PING is not compatible with WES data, or it was not tested in this manner.

kir-mapper is faster than PING but slower than T1K. Installation is easy and can be done directly on the system or using virtual environments. kir-mapper does not demand high processing power or large memory. It is possible to evaluate hundreds of samples simultaneously using a personal laptop with 16 Gb of memory. Like PING, kir-mapper supports analysing thousands of samples simultaneously and demands large sample sizes to define copy numbers accurately. However, for some genes, such as KIR2DL1, copy number definition is easier with PING than with kir-mapper. The genotyping tool of kir-mapper outputs all observed SNPs and InDels (known and new) for all samples in a VCF file using the hg38 genome as a reference, which can be embedded with the genotypes of the rest of the genome for WGAs. Unlike other methods, kir-mapper allows a manual inspection of the BAM files and the alignments for each gene. Therefore, when applying multiple software programs to detect KIR alleles (as we recommend), the user can inspect these BAM files in case

of inconsistent calls. Besides the BAM file, kir-mapper also produces text reports that allow the user to check the total count of tested and validated variants, the percentage of matches between the sample and tested alleles and a catalogue with the position of potential incompatibilities in the hg38 genome context. kir-mapper has a built-in method to solve ambiguities by using Shapeit4, which can be applied without any post-processing of the typing results. However, solving ambiguities works only when there are many samples. Therefore, like PING, kir-mapper is unsuitable for analysing a single sample or a small number of samples unless the user assumes that there are two copies of each KIR gene.

# 3.6 | Kir-Mapper Compatibility and Demands for Memory and Processing Power

We designed the kir-mapper to be compatible with personal computers, low-power workstations or high-power servers. It was written in C++, and it is compatible with Linux, macOS and Windows Subsystem for Linux (WSL). The minimum configuration depends on the number of samples to be processed simultaneously.

The exomes were processed using a late 2013 iMac with an i7 processor, 32 Gb of RAM and an HD of 3 Tb. WGS data was processed with an Ubuntu 22 workstation with a 12th-generation i9, 128 GB of RAM and 4 TB of SSD. We also genotyped the samples with KIR genes sequenced by the probe-capture method using a Windows laptop with a 12th-generation i7 processor, 16 Gb of RAM, 1 Tb of SSD and WSL2. The run time highly depends on the machine and the available processors, the amount of data to process and how fast the device can read and write data. kir-mapper is multithreaded and stores intermediate data in the disk to minimise memory use.

For instance, on the i7 iMac and the WES data, it realigned the reads in about 5 min per sample (the map function), called copy numbers for all samples in 4 min (the ncopy function), called SNPs and alleles for all KIR genes and samples in about 90 min (the genotype function) and the haplotypes to solve ambiguities in 2 min. The map function for WGS processed by the Ubuntu workstation took about 3 min per sample. The algorithm that is more demanding in terms of memory is the genotype (calling SNPs and InDels), and the amount of memory needed highly depends on the number of samples genotyped simultaneously and the number of threads used.

### 3.6.1 | Dependencies

kir-mapper depends on a series of third-party programs called by the main program when necessary. These include samtools, bcftools, BWA, freebayes, whatshap, shapeit4, picard tools, R and some R libraries to create the plots. The tested versions are listed on the program website. Installation with all dependencies can be done by using Conda/miniconda.

#### Acknowledgements

This work was supported by FAPESP/Brazil (Grant# 2021/14851-9) and the CAPES/COFECUB project (CAPES Project# 88881.879003/2023-01,

CAPES-COFECUB Me 1044/24). E.C.C. was supported by CNPq/Brazil (Grant# 307031/2022-5). P.J.N. was supported by NIH R01AI128775. R.N.P., G.S.P., M.R.F. and Í.S.F.S. are supported by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES, finance code #001) and CNPq/Brazil.

#### **Conflicts of Interest**

The authors declare no conflicts of interest.

#### **Data Availability Statement**

The data that supports the findings of this study are available in the Supporting Information of this article.

#### References

- 1. M. Colonna, A. Moretta, F. Vély, and E. Vivier, "A High-Resolution View of NK-Cell Receptors: Structure and Function," *Immunology Today* 21 (2000): 428–431.
- 2. N. K. Björkström, V. Béziat, F. Cichocki, et al., "CD8 T Cells Express Randomly Selected KIRs With Distinct Specificities Compared With NK Cells," *Blood* 120 (2012): 3455–3465.
- 3. C. Di Vito, J. Mikulak, and D. Mavilio, "On the Way to Become a Natural Killer Cell," *Frontiers in Immunology* 10 (2019): 1812.
- 4. K. L. O'Brien and D. K. Finlay, "Immunometabolism and Natural Killer Cell Responses," *Nature Reviews Immunology* 19 (2019): 282–290.
- 5. D. Pende, M. Falco, M. Vitale, et al., "Killer Ig-Like Receptors (KIRs): Their Role in NK Cell Modulation and Developments Leading to Their Clinical Exploitation," *Frontiers in Immunology* 10 (2019): 1179.
- 6. C. Fauriat, M. A. Ivarsson, H.-G. Ljunggren, K.-J. Malmberg, and J. Michaëlsson, "Education of Human Natural Killer Cells by Activating Killer Cell Immunoglobulin-Like Receptors," *Blood* 115 (2010): 1166–1174.
- 7. J. Dębska-Zielkowska, G. Moszkowska, M. Zieliński, et al., "KIR Receptors as Key Regulators of NK Cells Activity in Health and Disease," *Cells* 10 (2021): 1777.
- 8. A. M. Martin, E. M. Freitas, C. S. Witt, and F. T. Christiansen, "The Genomic Organization and Evolution of the Natural Killer Immunoglobulin-Like Receptor (KIR) Gene Cluster," *Immunogenetics* 51 (2000): 268–280.
- 9. A. M. Martin, J. K. Kulski, S. Gaudieri, et al., "Comparative Genomic Analysis, Diversity and Evolution of Two KIR Haplotypes A and B," *Gene* 335 (2004): 121–131.
- 10. M. Uhrberg, N. M. Valiante, B. P. Shum, et al., "Human Diversity in Killer Cell Inhibitory Receptor Genes," *Immunity* 7 (1997): 753–763.
- 11. P. J. Norman, J. A. Hollenbach, N. Nemat-Gorgani, et al., "Defining KIR and HLA Class I Genotypes at Highest Resolution via High-Throughput Sequencing," *American Journal of Human Genetics* 99 (2016): 375–391.
- 12. W. Jiang, C. Johnson, J. Jayaraman, et al., "Copy Number Variation Leads to Considerable Diversity for B but Not A Haplotypes of the Human KIR Genes Encoding NK Cell Receptors," *Genome Research* 22 (2012): 1845–1854.
- 13. J. Robinson, J. A. Halliwell, H. McWilliam, R. Lopez, and S. G. E. Marsh, "IPD—The Immuno Polymorphism Database," *Nucleic Acids Research* 41 (2013): D1234–D1240.
- 14. W. M. Marin, R. Dandekar, D. G. Augusto, et al., "High-Throughput Interpretation of Killer-Cell Immunoglobulin-Like Receptor Short-Read Sequencing Data With PING," *PLoS Computational Biology* 17 (2021): e1008904.
- 15. H. Li and R. Durbin, "Fast and Accurate Short Read Alignment With Burrows-Wheeler Transform," *Bioinformatics* 25 (2009): 1754–1760.

- 16. B. Langmead and S. L. Salzberg, "Fast Gapped-Read Alignment With Bowtie 2," *Nature Methods* 9 (2012): 357–359.
- 17. M. Byrska-Bishop, U. S. Evani, X. Zhao, et al., "High-Coverage Whole-Genome Sequencing of the Expanded 1000 Genomes Project Cohort Including 602 Trios," *Cell* 185 (2022): 3426.e19–3440.e19.
- 18. E. C. Castelli, M. A. Paz, A. S. Souza, J. Ramalho, and C. T. Mendes-Junior, "Hla-Mapper: An Application to Optimize the Mapping of HLA Sequences Produced by Massively Parallel Sequencing Procedures," *Human Immunology* 79 (2018): 678–684.
- 19. L. Song, G. Bai, X. S. Liu, B. Li, and H. Li, "T1K: Efficient and Accurate KIR and HLA Genotyping with Next-Generation Sequencing Data (2022), https://doi.org/10.1101/2022.10.26.513955.
- 20. G. F. Gao, D. Liu, X. Zhan, and B. Li, "Analysis of KIR Gene Variants in the Cancer Genome Atlas and UK Biobank Using KIRCLE," *BMC Biology* 20 (2022): 191, https://doi.org/10.1186/s12915-022-01392-2.
- 21. W. Marin and J. A. Hollenbach, Software Update: Interpreting Killer-Cell Immunoglobulin-Like Receptors From Whole Genome Sequence Data with PING. HLA1014414482023Software update: Interpreting Killer-Cell Immunoglobulin-Like Receptors from Whole Genome Sequence Data with PING 510.1111/tan.14949.
- 22. D. A. Benson, M. Cavanaugh, K. Clark, et al., "GenBank," *Nucleic Acids Research* 41 (2013): D36–D42.
- 23. P. Danecek, J. K. Bonfield, J. Liddle, et al., "Twelve Years of SAMtools and BCFtools," *GigaScience* 10, no. 2 (2021): giab008, https://doi.org/10.1093/gigascience/giab008.
- 24. M. S. Naslavsky, M. O. Scliar, G. L. Yamamoto, et al., "Whole-Genome Sequencing of 1,171 Elderly Admixed Individuals From São Paulo, Brazil," *Nature Communications* 13, no. 1 (2022): 1004, https://doi.org/10.1038/s41467-022-28648-3.
- 25. E. Garrison and G. Marth, "Haplotype-Based Variant Detection from Short-Read Sequencing," 2012, https://doi.org/10.48550/arXiv.1207.3907.
- 26. M. Martin, M. Patterson, S. Garg, et al., "WhatsHap: Fast and Accurate Read-Based Phasing," 2016, https://doi.org/10.1101/085050.
- 27. O. Delaneau, J.-F. Zagury, M. R. Robinson, J. L. Marchini, and E. T. Dermitzakis, "Accurate, Scalable and Integrative Haplotype Estimation," *Nature Communications* 10 (2019): 5436.
- 28. W. Huang, L. Li, J. R. Myers, and G. T. Marth, "ART: A Next-Generation Sequencing Read Simulator," *Bioinformatics* 28 (2012): 593–594.
- 29. T. D. J. Farias, S. Brugiapaglia, S. Croci, et al., "HLA-DPB1\*13:01 Associates With Enhanced, and KIR2DS4\*001 With Diminished Protection From Developing Severe COVID-19," *HLA* 103, no. 1 (2023): e15251.
- 30. W.-W. Liao, M. Asri, J. Ebler, et al., "A Draft Human Pangenome Reference," *Nature* 617, no. 7960 (2023): 312–324, https://doi.org/10.1038/s41586-023-05896-x.
- 31. T.-K. Hung, W. C. Liu, S. K. Lai, et al., "Genetic Diversity and Structural Complexity of the Killer-Cell Immunoglobulin-Like Receptor Gene Complex: A Comprehensive Analysis Using Human Pangenome Assemblies" (2023), https://doi.org/10.1101/2023.11.12.566753.
- 32. M. Stephens, N. J. Smith, and P. Donnelly, "A New Statistical Method for Haplotype Reconstruction From Population Data," *American Journal of Human Genetics* 68 (2001): 978–989.
- 33. F. F. Gonzalez-Galarza, A. McCabe, E. J. M. D. Santos, et al., "Allele Frequency Net Database (AFND) 2020 Update: Gold-Standard Data Classification, Open Access Genotype Data and New Query Tools," *Nucleic Acids Research* 48, no. D1 (2020): D783–D788, https://doi.org/10.1093/nar/gkz1029.
- 34. P. J. Norman, J. A. Hollenbach, N. Nemat-Gorgani, et al., "Co-Evolution of Human Leukocyte Antigen (HLA) Class I Ligands With Killer-Cell Immunoglobulin-Like Receptors (KIR) in a Genetically

Diverse Population of Sub-Saharan Africans," *PLoS Genetics* 9 (2013): e1003938.

#### **Supporting Information**

Additional supporting information can be found online in the Supporting Information section.