



Accuracy of NGS HLA typing data influenced by STR

Hanneke W.M. van Deutekom, Wietse Mulder, Erik H. Rozemuller*

GenDx, Yalelaan 48, 3584CM Utrecht, the Netherlands



ARTICLE INFO

Keywords:

STR
NGS
HLA typing
HLA-DRB5
HLA-DRB1
Repeat

ABSTRACT

Next Generation Sequencing (NGS) has become a major technology in HLA typing. The expectations are that highly accurate and unambiguous typing results will be obtained. However, HLA typing by NGS has some limitations caused by imperfections in the PCR amplification. The accuracy of NGS data is investigated by analyzing the Short Tandem Repeats (STR) regions. For this analysis HLA-DRB5 is used as the model. The repeat length in a sample highly influences the repeat length distribution present in the reads of NGS data. With a repeat length of 20 only 50% of all reads were of the estimated repeat length, seriously hampering distinguishing allelic differences in this region correctly. Our findings are confirmed by doing the same analysis in HLA-DRB1. Despite the uncertainty of determining the repeat lengths, several new HLA-DRB5 alleles have been identified in this paper.

1. Introduction

Next Generation Sequencing (NGS) has become a major technology in HLA typing. NGS is known for highly accurate and unambiguous HLA typing results. However, HLA typing by NGS has some limitations. Both the PCR amplification as well as the technique to determine the DNA sequence cause inaccuracies in the final reads. PCR errors may consist of randomly erroneous incorporation of single nucleotides and slippage, especially at homopolymer regions and regions with STRs [1–3]. Depending on the sequencing platform, additional errors may be encountered.

Since in NGS analysis each read originates from a single molecule, the aforementioned errors can be inspected in close detail. Using the detailed inspection the effect of those errors on the final analysis result is investigated.

To determine the influence of STR regions on the accuracy of the NGS data, the repeated region adjacent to the 3' end of exon two in HLA-DRB5 is studied. This region consists of two flanking dinucleotide repeats. The HLA-DRB5 gene is a good model because many individuals only possess a single copy of this gene, which eliminates analysis challenges caused by heterozygous individuals. HLA-DRB5 is associated with DRB1 * 15 or DRB1 * 16 alleles [4]. Selecting samples from individuals with either a DRB1 * 15 or a DRB1 * 16 allele, and a second allele which does not belong to these groups, guarantees that only one copy of the HLA-DRB5 gene is present. Based on the STR length distribution in the reads, an estimate of the repeat length can be determined. Additionally the percentage of reads with this estimated

repeat length can be identified. This percentage influences the reliability of the HLA typing, and could explain false-heterozygous calling in this region.

The analysis is extended to HLA-DRB1 gene, which contains a similar STR structure located in the beginning of intron two. As the repeat patterns within HLA-DRB1 are distinct, the alleles within the heterozygous samples can be separated.

2. Materials and methods

2.1. Sample selection

Samples were selected to contain a single copy of the HLA-DRB5 gene, or are homozygous for HLA-DRB1 and HLA-DRB5.

For DRB1, samples were selected to include samples of all DRB1 allele groups.

2.2. PCR

Human genomic DNA samples were used to amplify HLA-DRB5 (n = 18) and HLA-DRB1 (n = 23) by a Polymerase Chain Reaction (PCR) using NGSgo®-AmpX (GenDx) and the GenDx-LongRange PCR kit according to the Instructions for Use. The enzyme in this kit is a mixture of a Taq polymerase and a proofreading enzyme with 3'-5'-exonuclease activity for the removal of erroneously built-in, wrong bases. The processivity is primarily driven by the Taq polymerase.

The PCR consists of 35 cycles of the following steps: 15 s, 95 °C

* Corresponding author.

E-mail address: e.rozemuller@GenDx.com (E.H. Rozemuller).

<https://doi.org/10.1016/j.humimm.2019.03.007>

Received 31 May 2018; Received in revised form 6 March 2019; Accepted 6 March 2019

Available online 26 March 2019

0198-8859/ © 2019 Published by Elsevier Inc. on behalf of American Society for Histocompatibility and Immunogenetics.

denaturation, 30 s, 65 °C annealing, and 6 min, 68 °C elongation. To study the effect on PCR slippage, the number of PCR cycles in the amplification protocol increased in steps of 2 cycli from 25 to 35.

2.3. Library preparation and NGS analysis

HLA-DRB5 and HLA-DRB1 amplicons were processed in the NGSgo® workflow for Illumina using NGSgo-LibrX and NGSgo-IndX (GenDx). Libraries were paired-end sequenced (2x150bp) on a MiSeq platform (Illumina). FASTQ files were analyzed in NGSengine® HLA typing software V2.9 (GenDx), using IPD-IMGT/HLA database 3.31.0 [5] to determine the HLA-DRB5 typing based on exon 2. HLA-DRB1 typing was based exon 2, exon 3 and intron 3. The best match reported by NGSengine is the typing result which is best represented by the NGS data. Note that the repeat discussed in this paper is located in intron 2, and is thus excluded for obtaining the typing result.

2.4. Repeats

The DRB5 repeat can be represented as GGGAATCTGA (GT)_m (GA)_n GGAAGAGAGAG in which the factors *m* and *n* are numbers describing the lengths of the respective dinucleotide repeats.

A software tool was developed to analyze all reads from the NGS data which fully cover the repeat region, including the flanking sequences. For these reads, the *m* and *n* values are determined and the distribution of these factors is studied.

The estimated repeat length is defined as the value with the largest number of reads. This is based on the assumption that the majority of fragments that will be generated during the PCR reaction have the correct length. However, for longer repeat lengths, the majority of generated fragments have lengths which differ from the estimated repeat length and thus large numbers of fragments with an inaccurate length will be present (see Fig. 1B).

Data of HLA-DRB1 of 23 samples have been studied in a similar way. The DRB1 repeat is flanked by ATCTGA and CGCCAT. Within these two flanking regions there is a GT and a GA repeat, however the GA repeat is interrupted by other nucleotides. Hence, only the GT repeat has been analyzed, since this repeat is present in the same configuration in all alleles.

Note that, the GT repeat of DRB1 * 07:01:01 is interrupted by a G

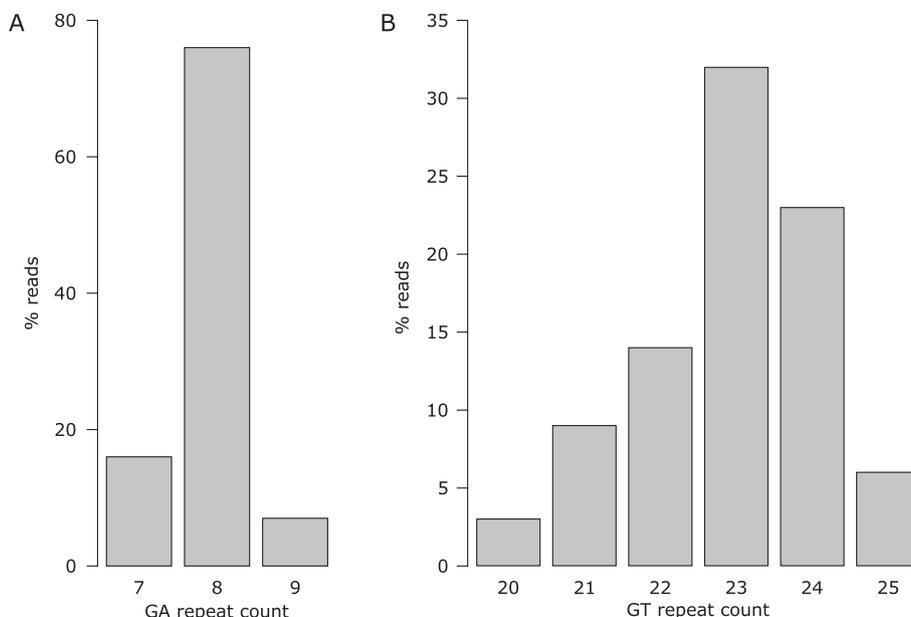


Fig. 1. Distribution of the length of the GA repeat and GT repeat in a sample typed as DRB5 * 01:02. 1A: The estimated GA repeat length is 8. 75% of all reads contained 8 times GA. 1B: The estimated GT repeat length is 23. 32% of all reads have 23 times GT.

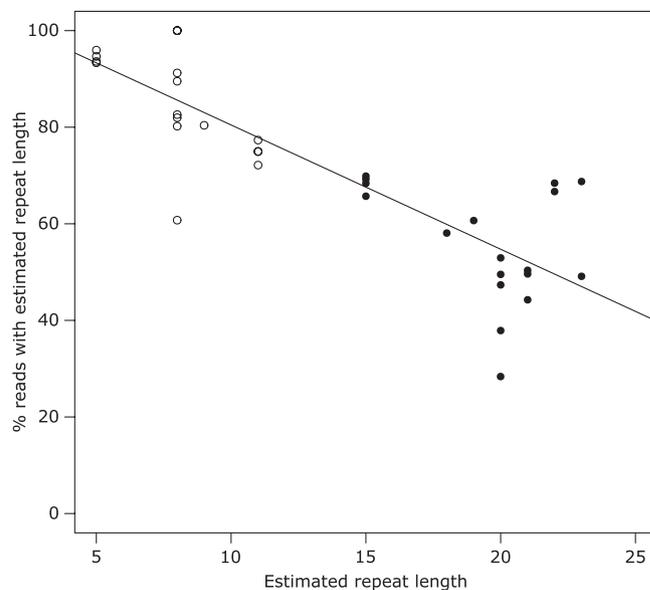


Fig. 2. Relation between the estimated repeat length and the percentage of reads with the estimated GA (open dots) and GT (closed dots) repeat length for HLA-DRB5.

and a T; ATCTGAGTGGTGT(GT)_m(GA)_nGCGAGACCGCCAT. The GT repeat used for this analysis starts after the additional T, as indicated in the pattern before. This pattern was observed in all samples with a DRB1 * 07:01:01.

18 of the DRB1 samples are heterozygous, due to the pattern around the GA repeat the two alleles of these samples could be distinguished reliably. For a list of the patterns of the repeats in DRB1 see Table 2.

3. Results

The NGS data of 18 DRB5 samples was used to determine the different lengths of the repeat per sample. A typical example of the length distribution of the GA repeat is shown in Fig. 1A. Here, 75% of all reads have a repeat length of 8 times GA, which is thus the estimated repeat length for this sample. The length distribution of the GT repeats in the

Table 1

The estimated repeat length for the DRB5 samples. For DRB5 * 01:01, four different configurations of the factors m and n are found. For DRB5 * 01:02 four different configurations are identified. For DRB5 * 02:02 a single configuration is found. The number of alleles identified for each combination is indicated between brackets. The * represents the configuration of the repeat in the genomic sequence of DRB5 * 01:01:01 within the IPD-IMGT/HLA database 3.34, which is the only repeat sequence in this database.

GA					
GT	m\n	5	8	9	11
	15				DRB5 * 02:02 (4)
	18		DRB5 * 01:01 (1)		
	19	DRB5 * 01:01 (1)			
	20		DRB5 * 01:01 (3)	DRB5 * 01:02 (1)	
	21	DRB5 * 01:01 (3)*	DRB5 * 01:02 (1)		
	22		DRB5 * 01:02 (2)		
	23		DRB5 * 01:02 (2)		

same sample is shown in Fig. 1B. The estimated repeat length is 23 repeats, which occurs in only 32% of all reads. 29% of all reads have 1 or 2 repeats more, 26% have less repeats.

For all 18 samples the relation between the estimated repeat lengths and the percentage of reads with this length is shown in Fig. 2. The estimated repeat length and the percentage containing this length are correlated ($p < 0.001$, $R^2 = 0.71$). Remarkably, with an estimated repeat length of 20, the percentage of reads with the estimated repeat length is only 50%.

Typing the samples considering only exon 2, i.e. using the core region, revealed 8 samples to be DRB5 * 01:01, 6 samples DRB5 * 01:02 and 4 samples DRB5 * 02:02. Analysis of the repeats in the 8 DRB5 * 01:01 samples showed 4 different repeat configurations, as is shown in Table 1. The number of occurrences of these configurations in the analyzed samples is indicated between brackets. Analysis of the repeats of the 6 DRB5 * 01:02 samples showed 4 repeat configurations. All DRB5 * 02:02 samples have the same repeat configuration of 15 times GT followed by 11 times GA.

In all samples the sequences flanking the repeats were identical to the sequences described in the Materials and Methods.

The analysis has also been applied to the HLA-DRB1 gene. However, the repeat in DRB1 is more complicated than that in DRB5. DRB5 has a clear $(GT)_m (GA)_n$ repeat, while the $(GA)_n$ repeat in DRB1 is interrupted by other nucleotides, creating patterns like “ATCTGA(GT)_mAA(GA)_xAA(GA)_yGCGCGCCAT”. A longer list of patterns is provided in Table 2. Because these patterns are different between allele groups, they can be used to separate the two alleles in a

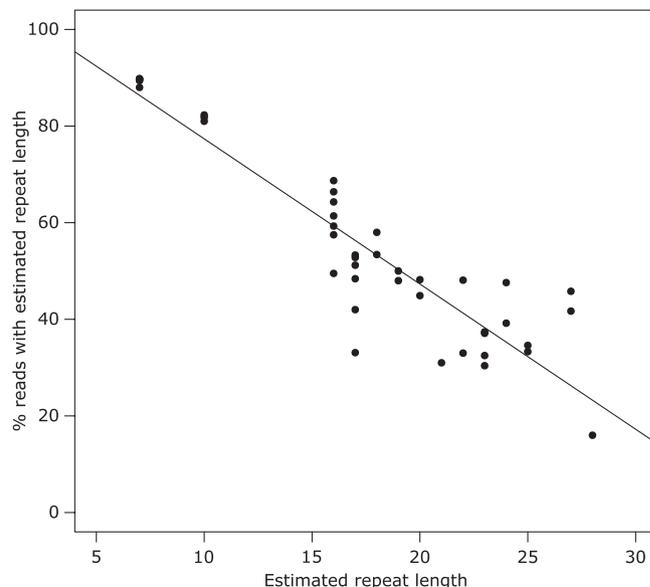


Fig. 3. Relation between the estimated GT repeat length in the HLA-DRB1 samples and the percentage of reads with the estimated repeat length.

heterozygous sample. Therefore, it is possible to determine the estimated GT repeat length in heterozygous DRB1 samples. Of the 23 samples that were used for analysis, 5 were homozygous and 18 were

Table 2

The DRB1 repeat patterns and estimated repeat lengths for the (3rd field) best matching alleles found in the panel. The sequence motifs flanking the repeats are underlined. Number of estimated repeat lengths are encoded by m, n1, n2 and n3.

Allele	Pattern	m	IPD-IMGT/HLA 3.34
DRB1 * 01:01:01	<u>ATCTGA</u> (GT) _m AAGAAA(GA) _{n1} GCGCGCCAT	16	16
DRB1 * 01:02:01	<u>ATCTGA</u> (GT) _m AAGAAA(GA) _{n1} GCGCGCCAT	19	15
DRB1 * 03:01:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CAGAGAGACA(GA) _{n2} GCGGCCAT	17/20	17/19/22
DRB1 * 03:02:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CAGAGAGACA(GA) _{n2} GCGGCCAT	19	19
DRB1 * 04:01:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} GCGGCCAT	22	21/22
DRB1 * 04:06:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} GCGGCCAT	22/23	22
DRB1 * 07:01:01	<u>ATCTGA</u> GTGGTGT(GT) _m (GA) _{n1} GCGAGACCGCCAT	7	7
DRB1 * 08:02:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CA(GA) _{n2} GCGCGCCAT	17	17
DRB1 * 08:03:02	<u>ATCTGA</u> (GT) _m (GA) _{n1} CA(GA) _{n2} GCGCGCCAT	17	17
DRB1 * 09:01:02	<u>ATCTGA</u> (GT) _m (GA) _{n1} GAAA(GA) _{n2} CAGAAAGAGGGAGCGCGCCAT	10	10
DRB1 * 10:01:01	<u>ATCTGA</u> CTCT(GT) _m (GA) _{n1} GCGCGCCAT	16/17	15/16/17
DRB1 * 11:01:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CAGAGAGACA(GA) _{n2} GCGGCCAT	23/24/25	25/27/28
DRB1 * 12:02:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CA(GA) _{n2} GCGGCCAT	28	28/30
DRB1 * 13:01:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CAGAGAGACA(GA) _{n2} GCGGCCAT	23	22/23
DRB1 * 13:02:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CAGAGAGACA(GA) _{n2} GCGGCCAT	20/21	20/24/26
DRB1 * 14:05:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CA(GA) _{n2} GCGGCCAT	25	26/28
DRB1 * 14:54:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CA(GA) _{n2} GCGGCCAT	24	21/25/34
DRB1 * 15:01:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} CA(GA) _{n2} CAGAGAGAGGAA(GA) _{n3} GCGGCCAT	18	17/18/19/20
DRB1 * 15:02:01	<u>ATCTGA</u> (GT) _m GAGACA(GA) _{n1} CAGAGAGAGGAA(GA) _{n2} GCGGCCAT	27	28/29
DRB1 * 16:01:01	<u>ATCTGA</u> (GT) _m (GA) _{n1} GAAA(GA) _{n2} GCGGCCAT	18	18

