

1 **PRESENTATION PAGE**

2

3 **ORIGINAL ARTICLE**

4

5 **TITLE**

6 shRNA design based on natural and drug induced variability targeting conserved  
7 regions from HIV reverse transcriptase

8

9 **RUNNING TITLE**

10 shRNAs directed against reverse transcriptase of HIV

11

12 **COMPLETE AUTHORS INFORMATION**

13 María Catalina Méndez Ortega<sup>1,2</sup>

14 Juan Carlos Mendoza Sánchez<sup>3</sup>

15 Luis Miguel Rodríguez Rojas<sup>2</sup>

16 Silvia Restrepo<sup>2</sup>

17 Gloria Janeth Rey Benito<sup>1</sup>

18

19 <sup>1</sup>Grupo de Virología SRNL del Instituto Nacional de Salud, Bogotá, Colombia

20 <sup>2</sup>Laboratorio de Micología y Fitopatología de la Universidad de los Andes, Bogotá,

21 Colombia

22 <sup>3</sup>VC@SOFT, Bogotá, Colombia

23

1 Correspondence should be addressed to M.C.M.

2 [mc.mendez159@uniandes.edu.co](mailto:mc.mendez159@uniandes.edu.co).

3

4 **Corresponding author :**

5 María Catalina Méndez Ortega, Grupo de Virología, Instituto Nacional de Salud,

6 Avenida Calle 26 No. 51-20, Bogotá, D.C., Laboratorio de Micología y

7 Fitopatología de la Universidad de los Andes (LAMFU)

8 Telephone: (1) 2207700 ext 439 - 341; fax number: (1) 2207700 ext 340

9 e-mail: mc.mendez@uniandes.edu.co

10

1

2 **ABSTRACT**

3

4

5 Inhibition of HIV-1 by RNA interference can be achieved by the expression of short  
6 hairpin RNAs (shRNAs) targeting conserved sequences. We designed shRNAs  
7 against HIV-1 reverse transcriptase considering the variability of conserved regions  
8 within the retrovirus reverse transcriptase conserved domain  $\alpha$ 01645, where the  
9 RNA dependent DNA polymerase activity of the enzyme resides. HXB2 (K03455)  
10 genome was used to choose regions with conserved residues important for  
11 enzyme activity. Regions were identified in multiple alignments from naive and  
12 drug-resistant isolates. A script was developed that predicted silencing efficiency  
13 based on previously reported parameters and that could identify highly frequent  
14 nucleotide combinations in a 21 base-long window within these regions. Higher  
15 number of sequence combinations was found in alignments from resistant isolates.  
16 shRNAs targeting reverse transcriptase active site residues W24 and P25 had  
17 scores over 7. Three-dimensional structural analyses from wild-type and resistant  
18 enzymes consistent with enormous variability explain the difficulties in finding a  
19 perfect region for RNAi mediated silencing. For clinical purposes, HIV-1 variability  
20 is an obstacle for efficient silencing from shRNAs designed towards a consensus  
21 sequence as there are too many functional variants of the enzyme. We suggest a  
22 complementary therapy using a combinatorial approach consisting of a mix of  
23 shRNAs targeting the most frequent viral variants from 1214 naïve genomes and  
24 1381 from resistant variants, during a highly active antiretroviral therapy regimen.

- 1 **Key Words:** Acquired Immunodeficiency Syndrome, HIV-1, RNA interference,
- 2 Gene Therapy, Antiretroviral Therapy, HAART
- 3

1

## 2 INTRODUCTION

3

4

5 Despite the advent of highly active antiretroviral therapy (HAART), human  
6 immunodeficiency virus (HIV-1) is still matter of public health concern [1, 2]. The  
7 major obstacle towards a cure lies in the integration of the viral genome [3-7].  
8 Consequently, the source of variability exists indefinitely and the virus will always  
9 have the chance to restart infection and become resistant leading to treatment  
10 failure and AIDS [8-13].

11 The overwhelming genetic variability in HIV-1 is principally due to the error-  
12 prone nature of reverse transcriptase (RT), characterized by its low fidelity with a  
13 calculated incorporation of one mutation per replication cycle [15, 22]. But other  
14 three mechanisms are also important: i) homologous and non-homologous  
15 recombination events between the two copies of the viral genome incorporated  
16 within the virion [14-17], ii) reverse transcriptase (RT) translocations during reverse  
17 transcription [18-21], and iii) high replication rate of the virus, enough to produce  
18  $10^{10}$  viral particles per day. These factors are responsible for the generation of  
19 quasispecies, that in conjunction with other forces contribute to the immune  
20 system exhaustion [6, 11, 23-26]. In addition, the virus source hides inside cellular  
21 reservoirs [8, 13, 27-29], which adds to its highly polymorphic nature, to its innate  
22 ability to accumulate mutations and to the flexibility of its proteins to accept them.  
23 Altogether these factors explain how the virus manages to overcome HAART [30,  
24 31]. Clearly, therapy or vaccine strategies should be efficient against each and

1 every possible competent-replication variant that may emerge even in presence of  
2 selective pressures different from the immune system [32-34]. Unfortunately,  
3 though HAART saves thousands of lives, emergence of resistant variants occurs,  
4 even though multiple key steps of the replication cycle of the virus are targeted  
5 simultaneously [9, 35-37].

6 RNA interference (RNAi) is an eukaryotic evolutionarily conserved natural  
7 occurring process by which double-stranded RNA (dsRNA) triggers post-  
8 transcriptional gene silencing [38]. Its existence is related to a defense mechanism  
9 against infectious or foreign nucleic acids (like viruses), to the control of mobile  
10 genetic elements, and to the regulation of gene expression [39].

11 RNAi is nowadays considered an appropriate alternative to HAART  
12 nonetheless; efficacy may be lowered in humans due to viral variability and  
13 resistance emergence. Several RNAi *in vitro* approaches have proven to be  
14 effective against HIV-1 [32, 40-45]. Furthermore, HIV-1 was silenced *in vivo* in  
15 humanized mice models [40, 46-49]. Certainly, a potent replication inhibition was  
16 obtained by the simultaneous targeting of multiple conserved regions in the viral  
17 genome with expression of multiple shRNAs [45, 50-53]. HIV variability underlies  
18 the fact that key target selection is of utmost importance. As an example, it was  
19 previously shown *in vitro* that when a sequence located within the *nef* gene, which  
20 is present in all the viral mRNA species, is targeted, a region including the  
21 sequence is deleted from the viral genome rendering the virus resistant to RNAi  
22 silencing [116]. This happens because that specific sequence is biologically  
23 unimportant, and this cannot be ignored even though the target was not thought to  
24 be Nef itself. Up to now shRNA design has been based on conserved sequences

1 from a reference sequence, on a multiple alignment and in some cases from  
2 previously in vitro screened silencing molecules from which the most efficient ones  
3 have been selected [50, 52, 54]. ter Brake et al 2006, found 19 highly conserved  
4 regions that were shared by 75% of 170 genomes [45] and in a similar approach  
5 Naito et al, found 216 sequences that showed conservation in more than 75% of a  
6 total of 495 nearly full-length genomes [55]. The most recent approach analyzed  
7 37,949 gene sequence fragments, from more than 8,000 genomes [56]. But no  
8 strategy has tried to understand and decipher HIV variability in terms of not trying  
9 to find a unique sequence for a specified conserved region, but instead in terms of  
10 identifying a region that has the least number of sequence combinations or viral  
11 variants from which the most frequent are chosen in order to design silencing  
12 molecules toward them. HIV variability makes impossible to have a unique  
13 sequence to target all the variants or quasispecies that can exist.

14 So even though silencing of HIV in humanized mice models undoubtedly  
15 showed the potential of RNAi for HIV-1 treatment, it is still far from being effective  
16 as unique therapy in an infected patient [46, 47, 49, 57-59]. HIV-1 mutates within  
17 the constraints of the protein structure-function relationship and a mouse life span  
18 is perhaps too short to evidence the natural outcome of resistant quasispecies.  
19 Resistance in humans against HAART inhibitors may arise within a few or more  
20 years, usually 2 to 4, of initiating therapy depending on different variables [8, 9, 60-  
21 62]. Then, longer periods of time are needed in an animal model in order to  
22 evaluate emergence of resistance or the inhibition of virus replication, even more  
23 when treatment has to be life-long. Unquestionably, shRNAs must be carefully  
24 selected [63].

1 We propose a shRNA design approach considering all the variability from HIV  
2 sequences available in non-redundant HIV databases. We aimed to find a  
3 conserved region in which few shRNAs can be used to target all the sequences in  
4 a multiple alignment, using a Score-based attempt in order to aid selecting for  
5 efficient molecules. Our approach is based on identifying all the nucleotide  
6 changes that may be accepted in conserved 21-nucleotide long regions localized  
7 within conserved domains of HIV-1 RT. For this, we aimed to detect changes  
8 occurring during selective drug pressure by analyzing multiple sequence  
9 alignments performed with resistant isolates sequences retrieved from HIV Drug  
10 Resistance Database [64-66], based on a first line RT inhibitor (RTI) treatment. A  
11 curated multiple sequence alignment from HIV databases was also used to search  
12 for natural occurring changes. As nucleotide changes cannot accumulate in an  
13 infinite fashion and have to follow structure restrictions to preserve enzyme  
14 function we hypothesized that a set of shRNAs might exist that can silence specific  
15 virus variants and possibly all of them. Our main target is the RT conserved motif  
16 YMDD [67] that lies within the retrovirus reverse transcriptase conserved domain  
17 cd01645 RT\_Rtv [67, 68] which is responsible for the RNA dependent DNA  
18 polymerase activity, required for the viral replication cycle [67, 69] and thus we  
19 expect the least number of changes to be allowed.

20

21

22

23

24



1

## 2 **MATERIALS AND METHODS**

3

4

### 5 **Identification of conserved regions and possible target sites within the**

#### 6 **RT\_Rtv domain**

7 The search for conserved domains was performed in Pfam and InterProScan  
8 databases [70-74]. Identified domains were then localized in the reference  
9 sequence HXB2 accession number K03455  
10 (gi|1906382|gb|K03455.1|HIVHXB2CG Human immunodeficiency virus type 1  
11 (HXB2), complete genome; HIV1/HTLV-III/LAV reference genome). This sequence  
12 was used as query against the conserved domain database of NCBI [75], in order  
13 to identify catalytic residues, residues involved in DNA binding, dNTP binding, or  
14 active site residues without any other annotation. The map of coordinates from  
15 HXB2 reference sequence available at HIV databases was used to find the exact  
16 position of selected regions in the alignments, using the sequence as a guide  
17 (<http://www.hiv.lanl.gov/content/sequence/HIV/REVIEWS/HXB2.html>). Overview of  
18 the multiple alignments was also used to identify other conserved regions.

19

20

21

#### 22 **Identification of study regions in Resistant Isolates**

23 In addition to the mentioned regions from the conserved domain cd0548, high  
24 prevalence resistance mutations were selected from the "Mutation prevalence

1 According to Subtype” <http://hivdb.stanford.edu/cgi-bin/MutPrevBySubtypeRx.cgi>  
2 web page, from the Stanford HIV Drug Resistance Database [75, 76]. Mutations  
3 that localized within the conserved domain selected regions were preferred, but  
4 other mutations with high prevalence (>10% from the 3583 sequences used) were  
5 also chosen for analysis, if they mapped as consecutive amino acids. The following  
6 parameters were used to retrieve the data: mutation threshold  $\geq 1$  % AND  $\geq 2$   
7 persons, no mixtures (to exclude possible sequencing errors), and position  
8 mutation count = one mutation per person (more than one isolate per person if  
9 different mutations are observed at the same position of the isolates). The same  
10 exact parameters were used for NRTI-experienced vs RTI naive, NNRTI-  
11 experienced (+/-NRTI) vs RTI naive, and RTI-experienced vs RTI naive.

12

13

14

### 15 **3D Structure Analysis**

16 Thirty crystallographic structures of HIV-1 RT available at [www.rcsb.org](http://www.rcsb.org) protein  
17 databank (PDB) [77, 78] were used to perform Structural Alignments for  
18 comparison between different mutant enzymes, and wild type enzymes bound to  
19 different inhibitors. Heteroatoms were not included in the structural alignment. The  
20 following crystallographic structures were used: wild type structures 2ZD1 (1.8 Å)  
21 [79], 2RF2 (2.4 Å) [80], 3DLE (2.5 Å) [81], 2OPP (2.55 Å) [82], 1TKT (2.6 Å) [83],  
22 1TKZ (2.81 Å) [83], 1TKX (2.85 Å) [84], 1TL1 (2.9 Å) [83], 1TL3 (2.8 Å) [83], 1S1X  
23 (2.8 Å) [85], and mutant structures 3BGR (2.1 Å) [79], 3DLK, (1.85 Å) [86], 1S1T  
24 (2.4 Å) [86], 3DOL (2.5 Å) [81], 3DMJ (2.6 Å) [81], 1S9E (2.6 Å) [87], 1S1V (2.6 Å)

1 [85], 1S1W (2.7 Å) [85], 2OPQ (2.8 Å) [82], 1FKO (2.9 Å) [88], 1FKP (2.9 Å) [88],  
2 3DOK (2.9 Å) [81], 2OPR (2.9 Å) [82], 2ZE2 (2.9 Å) [79], 1S6P (2.9 Å)[87], 1S1U  
3 (3.0 Å) [85], 2IC3 (3.0 Å) [89], 3DM2 (3.1 Å) [81]. Protein structural alignments  
4 were generated with MultiProt [90] and analyzed in PDBviewer (Deep View) [91,  
5 92]. Mutant (MUT) and wild type (WT) structural alignments were performed  
6 independently and together. In either case 2ZD1 structure was used as reference  
7 layer. Protein structural alignments generated with MultiProt were run online,  
8 introducing PDB structure codes. Structure PDB codes were chosen to be kept in  
9 order to ensure that the first protein in the pdb file (the reference layer) was 2ZD1.  
10 A 2.5 Å threshold was chosen to generate the structural alignment with MultiProt.  
11 The comparison between WT and MUT structural alignments was performed with  
12 DeepView.

13

14

## 15 **Sequence Retrieval**

16 Curated sequences were retrieved from <http://hivdb.stanford.edu/> and from  
17 <http://www.hiv.lanl.gov/content/index>. Resistant isolates sequences were retrieved  
18 from the Stanford HIV Resistance Database [66, 93, 94] by a therapy criterion that  
19 consisted of different combinations of three antiretrovirals, two nucleoside analogs  
20 reverse transcriptase inhibitors (NRTI) and one non-nucleoside analog reverse  
21 transcriptase inhibitor (NNRTI) respectively. Duplicates and repeated sequences  
22 were eliminated using ElimDupes [http://www.hiv.lanl.gov/content](http://www.hiv.lanl.gov/content/sequence/ELIMDUPES/elimdupes.html)  
23 [/sequence/ELIMDUPES/elimdupes.html](http://www.hiv.lanl.gov/content/sequence/ELIMDUPES/elimdupes.html) tool available at HIV databases  
24 [<http://www.hiv.lanl.gov/>]. Sequences were further processed to eliminate foreign

1 characters not compatible with FASTA format (e. p. “~”). A total of 1250 sequences  
2 of HIV-1 isolates from HIV databases curated alignment were included in the study,  
3 as well as 2263 sequences from the Stanford University HIV-1 resistance database  
4 [64-66, 93, 94]. Multiple alignment was performed with HIV align tool from HIV  
5 Databases, using Mafft [95].

6 In addition, two first line regimens were selected from the National Resolution No.  
7 3442-2006 for the care of the HIV infected patient, approved for HIV infection  
8 treatment in Colombia by the Ministerio de Proteccion Social (Resolution 3442 for  
9 the Care of the HIV patient) [96] in order to retrieve sequences from resistant  
10 isolates generated during each of these drug selective pressures. The regimens  
11 are zidovudine-lamivudine-efavirenz (ZDV-3TC-EFV), which is normally the first  
12 choice, and stavudine-lamivudine-nevirapine (D4T-3TC-NVP), which is also a first  
13 line regimen.

14

15

## 16 **Colombian Resistant Isolates**

17 Twenty isolates from symptomatic HIV positive patients were included in the study,  
18 and the fasta sequences were incorporated in the group of resistant viral  
19 sequences. Only samples with informed consent, complete clinical history and viral  
20 loads over 1000 copies/ml were included. Virions from 1ml of plasma were  
21 concentrated at 15000 rpm x 99 min, resuspended in 500 ul of lysis buffer  
22 (Guanidinium thiocyanate - Biomerieux) for 10 min at room temperature. RNA  
23 extraction was performed using the semi automated Nuclisens Minimag, and  
24 Nuclisens extraction Kit, following manufacturer's instructions (Biomerieux, Marcy

1 l'Etoile, France) and eluted in 30 ul, of elution buffer. HIV genotyping was  
2 performed using TRUGENE HIV-1 commercial kit (Siemens Diagnostics, Deerfield,  
3 USA), and fasta sequences were further analyzed with HIV Stanford Drug  
4 Resistance Database software [97]. Quality assessment and CPG were used to  
5 evaluate probable sequencing errors and mutation prevalence respectively [75].

6  
7

### 8 **Multiple Alignments**

9 Curated multiple alignments are premade alignments already edited by HIV  
10 experts and are freely available at HIV databases web site  
11 <http://www.hiv.lanl.gov/content/sequence/NEWALIGN/align.html>. Not curated  
12 multiple alignments from resistant isolates were performed with the HIVAlign tool  
13 from HIV databases, which is an implementation of the HMMER package  
14 <http://www.hiv.lanl.gov/content/sequence/HMM/HmmAlign.html> . Sequences were  
15 codon aligned with GeneCutter [[http://www.hiv.lanl.gov/content/sequence/GENE\\_](http://www.hiv.lanl.gov/content/sequence/GENE_CUTTER/cutter.html)  
16 [CUTTER/cutter.html](http://www.hiv.lanl.gov/content/sequence/GENE_CUTTER/cutter.html)], and reference sequence HXB2 was included in the  
17 alignment as guide sequence to help find selected regions. Selected regions were  
18 found using alignment editor Bio-edit and eBIOX. Four multiple alignments were  
19 used for the study, one downloaded premade-curated alignment and three  
20 generated in this study (resistant isolates). Pre-made alignment corresponds to  
21 HIV-1 Group M plus recombinants whole web 2008 sequence alignment (not  
22 resistant isolates), from HIV Databases; not pre-made alignments correspond to  
23 resistant sequences alignments.

24

1

2 **Identification of nucleotide changes for each position of the 21 nucleotide**  
3 **long selected target region**

4 A Perl script was designed to identify changes in each position of the target  
5 sequence for all the sequences aligned along each multiple alignment or the whole  
6 multiple alignment (wide and large) using a window size of 21 nucleotides. The  
7 script could find unique combinations within the selected region of the alignment  
8 and their frequencies, showing the nucleotide changes. The script was designed to  
9 include the ambiguous bases from the 21-nucleotide window interpreting each of  
10 them into all the possible combinations according to the IUPAC code and finally  
11 adding them to the frequency of its corresponding unique combination.

12

13

14 **Design of the shRNAs that include all the variability of HIV-1 in a 21-**  
15 **nucleotide long target region**

16 The same script included a Score-based program that included parameters that  
17 have been frequently found in efficient shRNAs and siRNAs and suggested by  
18 Reynolds [98] and some by Jagla et al. 2005 [99] , but in this case, the efficiency  
19 information was searched and found in the DNA sequence. That is, the program  
20 gave a score of how efficient will be the silencing induced from the RNA  
21 transcribed from the DNA sequence (template guide strand) or unique combination  
22 of the target sequence. Furthermore, the program showed which parameters were  
23 not fulfilled and which did. A score of 11 points meant that all the parameters were  
24 successfully accomplished. The script followed 6 parameters which were: wins

1 (window size), star\_pos (position at which the script begins), end\_pos (position at  
2 which script ends), threshold (tells to discard the window if the number of  
3 combinations is higher than the chosen percentage of sequences), scorethres  
4 (minimum score allowed in the window; score is based in Reynolds conditions [98],  
5 min\_freq (minimum frequency allowed for each unique combination),  
6 min\_freq\_count (tells not to count the window if the chosen min\_freq\_count  
7 percentage of sequences do not meet min\_freq). Percentages were selected  
8 based on the number of sequences present in each multiple alignment.

9  
10

11 **RESULTS**

12  
13

14 **1. Identification of conserved regions and target sites within the**  
15 **conserved domain RT\_Rtv**

16 **Conserved Domains and YMDD motif**

17 Search was performed using Pol poliprotein AAB50259.1 from HXB2 reference  
18 sequence as query, and four conserved domains were found. Conserved  
19 domain RT\_Rtv was identified with an E-value of 1e-85, in CDD NCBI [68, 100].  
20 Within this domain 9 subregions were mapped that grouped DNA binding sites  
21 (green), dNTPs binding sites (dark blue), reverse transcriptase inhibitors (RTIs),  
22 binding sites (cyan blue), active sites with no other annotations (purple) and  
23 important motif YMDD (pink) [101] (Figure 1a). Table 1 shows the regions  
24 chosen for the analysis, their position in HXB2 reference genome, the important

1 residues within them including their position with respect to RT enzyme (HXB2  
2 numbering system), and the nucleotide sequence for each region.

## 3 4 **2. Resistance analysis from Colombian isolates.**

5 A subset of 50 samples from 130 hospitalized symptomatic HIV positive  
6 patients, collected from January of 2008 to October of 2009 were analyzed in  
7 order to include the sequences from resistant isolates in the multiple alignment  
8 from resistant sequences. From those 50 sequences, only samples with viral  
9 loads over 1000 copies/milliliter were chosen for genotyping, for a total of 30.  
10 One sequence went out of study (the patient leaved the study voluntarily), 9  
11 were repeated, and three were not recognized by COBAS AMPLICOR HIV-1  
12 MONITOR V1.5 (ROCHE) (no viral load data for these samples) but were  
13 successfully genotyped. In total, only 25 from the 30 samples resulted in  
14 TRUGENE HIV-1 report (Siemens Diagnostics, Deerfield, USA). One sample  
15 was diagnosed as having multiple quality problems (e. i. stop codons - frame  
16 shifts – too many unusual residues), the sequencing run was revised and  
17 complete assay repetition was not enough to achieve a higher quality and was  
18 not included (Sup Table 1). Thirty six percent of the protease (PR) sequences  
19 had any mutation, 12% of the RT sequences had any protease inhibitor (PI)  
20 mutation, 24% had any nucleoside reverse transcriptase inhibitor (NRTI)  
21 mutation, 16% had any NNRTI mutation, 12% had any of both types of  
22 mutations. Four percent of the PRRT sequences had any NRTI, NNRTI and PI  
23 3-class resistance mutations [102]. PI mutations seen in the data set were I54T  
24 and M64I, with a set prevalence of 4% and 8%. Major NRTI thymidine analog



1 mutations (TAMs) were D67N (4%), M184V (16%), L210W (4%), and T215Y  
2 (4%), major non-thymidine analog mutations were K65R (8%) and K70R (4%),  
3 and F77L (4%) was the only multi-NRTI resistance mutation seen in the data  
4 set. Major NNRTI mutations were K103N (8%), Y181C (4%), Y188L (4%), and  
5 P225H (4%). Only one mutation G152R, indicative of APOBEC3G –mediated  
6 editing was seen, in only one sequence. No more hyper mutated sites were  
7 identified, and at least one atypical mutation was observed for each of the  
8 sequences analyzed. Identity of PR sequences with HXB2 reference sequence  
9 was of 95.16%, for RT of 94.29%, and for the complete sequences PRRT of  
10 78.68% (mean value). Five samples showed high level resistance to NRTIs,  
11 four of these against 2 NRTI and one against 3 NRTIs. No high level resistance  
12 against PIs was observed for the studied set of samples (Sup Table 2).

### 14 **3. Resistance mutations within selected regions**

15 Resistant mutations localized near or within the selected regions were identified  
16 in the corresponding alignments. Table 2 shows the wild type residues, their  
17 position with respect to reference sequence HXB2, the mutated residue and  
18 prevalence based on HIV Drug Resistance Database data. Regions are named  
19 based on table 1. Our hypothesis is based on that if mutations are selected by  
20 certain regimen they could become “fixated” under selective drug pressure, and  
21 in the fact that if they match in mutable positions exactly within conserved  
22 active site residues, this could increase the “conservation” of the whole site,  
23 making it a possible target for gene therapy. From the regions that have  
24 mutations within or next to them, the most appropriate to target are those which

1 have the same mutation selected by both types of antiretrovirals (NRTIs –  
2 nucleoside RT inhibitors, NNRTIs - non nucleoside RT inhibitors), and even  
3 better if there is only one observed mutation for the residue. Following this  
4 criterion, chosen regions were number 4, 7, and 8. Since the targeted  
5 nucleotide region is about 21 nts long (7 triplets), regions 4 and 7 must include  
6 an extra residue at any end. Nonetheless, the other regions were screened and  
7 in case of region 9, since the residues next to each end of the initially chosen  
8 residues are not conserved and are mutable under drug selective pressure,  
9 windows including both residues were analyzed.

#### 11 **4. 3D structure analysis**

12 Structural alignment showed variability between mutant and wild type RT  
13 enzymes. Thirty crystallographic structures were used for the alignments.  
14 Structures had a mean residue number of 942.55, which ranged from 914 to  
15 979 residues. Structure Resolution ranged from 1.8 Å to 3.1 Å. Structure 2ZD1,  
16 was used as reference layer since it has the higher resolution 1.8 Å, and  
17 corresponds to a wild type enzyme. Figure 1b shows the structural alignment of  
18 wild type RT enzymes, with 774 aligned residues and 0.81 root mean square  
19 deviation (rms). Structural alignment evidenced high variability between non-  
20 mutant RTs. Being structures of the same wild type enzyme from the same  
21 virus we expected a better alignment, also that structures were not so variable  
22 between them and thus a bigger number of residues participating in the  
23 alignment. Figure 1c corresponds to the structural alignment of resistant  
24 (mutant) enzymes that showed even higher variability, with only 590 residues

1 participating in the alignment and 0.88 ms. This variability is evident in the  
2 superimposed alignments showed in Figure 1d. Structures have different  
3 mutations, and are also in complex with different ligands. Variability is evident in  
4 the number of residues involved in the alignment, which is far less than the  
5 number used for wild type RT structure alignment.

6 Even though structural alignments were performed with Deep View and  
7 Multiprot, WT and MUT structural alignment corresponds to the best generated  
8 with Multiprot (best from three).

## 10 **5. Multiple Alignments**

11 A total of 2533 sequences were retrieved from the Stanford drug resistance  
12 database, 286 were duplicates, and thus were eliminated. After editing  
13 characters, multiple alignment was generated with HIV align tool from HIV  
14 Databases tools  
15 (<http://www.hiv.lanl.gov/content/sequence/HMM/HmmAlign.html>) which is an  
16 implementation of HMMER package. Sequences were aligned by HMM (Hidden  
17 Markov Models) model HIV-1/SICcpz, allowing 5 codons to compensate for  
18 frame shift [103]. In addition, sequences selected under the regimens ZDV-  
19 3TC-EFV and D4T-3TC-NVP were also retrieved from this database.

20 Sequences for the premade curated alignment were retrieved from the HIV  
21 Databases, and consisted of curated sequences updated to 2008 and  
22 correspond to Group M HIV-1, including recombinant forms. This multiple  
23 alignment was not further edited.

## 6. shRNAs guide strands that include the variability found.

### Curated Alignment

Along the curated multiple alignment from Group M +Recombinants 2008 (HIV databases), only 2 regions arose a combination of Scores and number of unique combinations (frequency) that could be useful for gene therapy (e. i. 1 and 7). Different conditions of the script parameters were tested and found that there was no region without, at least one unique combination with a score of 4. A higher threshold for minimum score rendered no results. From all the different conditions evaluated, the following were finally used to test the selected regions from curated alignment: threshold 20%, Score-threshold = 4, min\_fre\_count 10%, min\_freq = 6. Results had to be exhaustively examined because not all the windows thrown were useful.

Table 3 depicts the results for region 3230- 3243, windows 3234 and 3237, and for region 3648 – 3662, window 3648. In all the alignment, there were no other regions with less than 46 unique combinations in which the most frequent ones rendered minimum scores of 7, the minimum score required for silencing.

Region 3230-3243, window 3234, with 46 total unique combinations, had scores between 4 and 9, but the most frequent unique combination occurred 959 times, with a score of 8. The next most frequent unique combinations occurred 40, 32, 29 and 24 times in the alignment, and only one had a score under 7, the one that occurred 40 times. Together without taking into account this last combination, the group of sequences accounted for the 89.3 % of the 1256 sequences in the alignment. All the remaining combinations had frequencies under 10, mostly of 1 (this is also true for the other regions and

1 windows obtained). About the score conditions, none of the most frequent  
2 combinations met conditions 1, 6 and 8, that correspond to having an A in  
3 position 3, and having an A or T in positions 15 and 17.

4 Window 3237, with 54 unique combinations showed better results, being only  
5 between 8 and 9, having 9 for the most frequent combination, which occurs 942  
6 times. The next most frequent unique combinations occurred 40, 32, 27, and 23  
7 times, all of them with a score of 8. All together these combinations accounted  
8 for the 87.6 % of the sequences in the alignment. In this case, score conditions  
9 1 and 11, having A in position 3 and having A or T in position 20, were the only  
10 conditions not met by all the 5 most frequent unique combinations. Individually,  
11 combinations with a score of 8 had one extra condition which they did not meet.  
12 Window 3548, from region 7, with 46 unique combinations, had the best scores,  
13 being between 9 and 10, but had twice the number of most frequent unique  
14 combinations with respect to the other windows. The frequency of the most  
15 common unique combinations was more equivalent between each other, being  
16 167, 161, 27 and 25 for scores of 10, and 255, 165, 113, 91, 25 and 22, for  
17 scores of 9. All together correspond to the 86.6 % of the sequences.

18 Other three windows from region 3391 – 3420 gave a number of unique  
19 combinations under 61, and scores of 7, 8, and 9 but included a gap. The same  
20 was true for region 3628- 3662, window 3648, which scored over 8. Gap editing  
21 did not improve the results probably because nucleotide positions were  
22 changed, changing the whole score.

23  
24 **Resistant sequences**

1 The script threshold conditions used for analyzing the curated alignment were  
2 not useful for analyzing the resistant sequences multiple alignment. Using these  
3 same conditions the analysis rendered too many windows. The only condition  
4 kept was score threshold at 4. In fact conditions had to be adequated to each  
5 alignment since each one has a different number of sequences and in turn  
6 representative variability. For each regimen alignment it was possible to find  
7 target regions with fine scores, specifically in the same regions in which it was  
8 possible to find targets for the curated alignment. We find the highest score of  
9 11, and also of 10, but rarely frequently. The small number of combinations in  
10 target sequences for shRNA design found for DAT-3TC-NVP seems to be a  
11 matter of number of sequences present in the alignment. Nonetheless the  
12 number of sequences in the ZDV-3TC-EFV is close to the one of the curated  
13 alignment, and gave more combinations than the alignment of the other  
14 regimen, even when the curated alignment has sequences from all group M  
15 plus recombinants (Table 3).

16 Curated alignment is supposed to be more divergent than regimen  
17 sequence alignments, which are principally, but not uniquely composed of B  
18 subtype strains. In the case of these alignments, silencing would be determined  
19 by the number and selection of the target sequences that can be introduced  
20 into a gene therapy vector. For D4T-3TC-NVP the number of combinations  
21 accounted for 78.02% (window 179), 87.9% (window 264), 83.52% (window  
22 459), 73% (window 565) and for ZDV-3TC-EFV alignment the number of  
23 combinations accounted for 89.07% (window 72) of the variability present in the

1 alignment, and 88.06% for window 73. Results for each window may seem  
2 different, but only *in vitro* studies will open light in the true silencing efficiency.

## 4 **7. BLAST**

5 Blastn searches were performed from CLC DNA Workbench Desktop for each  
6 of the unique combinations in each region, against Human genomic plus  
7 transcripts (test/gpipe/9606) all contig\_and\_rna NCBI database (table 4), and  
8 against Human subset of ESTs database (table 5). Search was limited by  
9 Entrez Organism *Homo sapiens*, using low complexity and mask for lower case  
10 as filters. In addition, default values for word size, match, mismatch, gap open  
11 cost, and gap extension cost, were used. Only the minus strand hits found are  
12 shown. In all cases Hits mapped to unique sequence tagged sites (in human  
13 chromosomes, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20,  
14 X genomic contigs, reference or alternate assembly), in no case an overlap of  
15 100% was seen. Most of the hits repeated along the region, but most E-values  
16 were not significant (1.2 -4.5). Only most significant E-values of 0.08 and 0.3  
17 are shown (Tables 4 and 5). Hits on minus strands are important because  
18 expressed sequences from these hits could be targeted. From the significant E-  
19 values (0.3), 11 targeted minus strands, and from the two most significant E-  
20 values (0.08), both targeted minus strands. In the case of targets from resistant  
21 sequences, number of hits ranged from 5 to 100 for each combination, each  
22 having at least one significant hit with significant E-value MISMA TABLA. Only  
23 important ESTs hits are shown for AZV-3TC-EFV, like IL-2, white blood cells, B  
24 cells etc., since many of the hits were ESTs from different types of cancer in

1 different tissues not infected by HIV. Important hits are those expressed by  
2 gene therapy target cells or progenitors.

#### 4 **DISCUSSION**

5 This is the first approach to a novel design of shRNAs guide strands, based  
6 on the scored search of a group of sequences that target most or all the subset of  
7 wild-type and drug-resistant RT variants. We developed a script in order to follow  
8 the established parameters of effective guide strand for shRNAs from Reynolds et  
9 al.2004 [98]. In our approach shRNAs guide strand sequences are directed to  
10 silence resistant variants selected during antiretroviral therapy, avoiding therapy  
11 failure and consequences that result from it. A set of guide strands was also find  
12 for HIV-1 Group M plus recombinant HIV-1 viruses.

13 shRNAs design is difficult due to the multiple requirements to achieve an  
14 efficient *in vivo* silencing and to all the parameters that had to be carefully followed.  
15 Available programs are normally directed to siRNA and not shRNA design [104-  
16 106], and some authors have shown that these are not useful to correctly predict  
17 the silencing efficiency of shRNAs, perhaps because processing occurs through  
18 different cellular pathways [107-109]. Developed tools like Gene Link  
19 (<http://www.genelink.com/sirna/shRNAi.asp>) do not allow the use of more than one  
20 aligned sequence, which is a necessary issue when trying to design silencing  
21 molecules against error-prone viruses like HIV. Recently miRanda [110] was used  
22 to check *in silico*, target sites for miRNAs in the 3' UTR of the HIV-1 molecular  
23 clone LAI genome, which were then tested with rna22 [111, 112]. This approach  
24 was based on using miRNAs and shRNAs with partial complementarities instead of



1 a total identity, which was thought to be absolutely necessary for shRNAs efficient  
2 silencing. Because of this premise, there was no need for a multiple alignment but  
3 it surely would have been informative.

4 The search for highly conserved regions in HIV-1 genome has been done  
5 for several RNAi screening, and has been based on multiple alignments and other  
6 types of sequence comparisons [45, 51, 56]. Nonetheless, no approach has  
7 analyzed in detail the conservation based on protein function, which is thought to  
8 be the least variable. In multiple alignments of HIV nucleotide sequences the  
9 observed variability is a result of the number of sequences present in the  
10 alignment, unless the whole set of sequences analyzed are representative of the  
11 virus variability which will be more than a million sequences for a conserved region  
12 based on the different number of DNA codons for each specific amino acid. In the  
13 particular case of HIV-1, multiple alignments from HAART treated isolates as well  
14 as from untreated isolates reveal enormous variability even on very conserved  
15 regions. As it was evident in the multiple alignments from this study, variability in  
16 the most restricted or conserved regions is still high enough to make any design  
17 strategy based on finding a consensus shRNA guide strand sequence inefficient in  
18 terms of long term or indefinite silencing, especially in a human HIV-1 infected  
19 patient. More over, even though multiple shRNAs targeted to conserved regions of  
20 the viral genome have proven to inhibit HIV-1 replication, viral escape has shown  
21 to be unimpeded [113-117]. Hypothetically, in an infected patient all the possible  
22 variants of each of the targeted conserved sequences must also be targeted in  
23 order to achieve life-long silencing. Eventually, this would increase the number of  
24 shRNAs that would have to be expressed from a single or multiple vectors, which

1 may be a limitation if the number and size exceeds certain limit [118]. The targeting  
2 of multiple conserved regions remains a constant strategy since it has showed  
3 encouraging results [43, 45, 51, 53, 54, 56, 118, 119], even though one should ask  
4 if silencing efficiency is the result of targeting one same sequence twice or more  
5 (times the number as conserved sequences are being targeted), due to targeting of  
6 more number of viral variants once with the only conserved sequence guide strand  
7 from all to which complementarity occurred, or both.<sup>3</sup>

8           Still there is another question: if all escape variants are indeed the result of  
9 “resistance” to RNAi therapy or if they are naturally occurring variants which are  
10 not being targeted and instead are being selected.

11           To counteract this problem, we decided to target a unique highly conserved  
12 region that tested with our bioinformatics tool showed the lowest number of  
13 different forms or unique combinations along the multiple alignments, that met a  
14 score of at least 7 for the most frequent ones, and that this most frequent variants  
15 were very abundant so that all together represented the 80% or more of the  
16 sequences present in the alignments. Conserved regions were chosen based on a  
17 protein-function approach, aiming to find a region with near seven consecutive  
18 conserved amino acids. Seven is the minimum score required for silencing [98].  
19 Unfortunately, the highest score (11) was not frequent, and even worse, it was  
20 mainly found in low conserved regions (data not shown). Consequently shRNAs  
21 requirements were easily accomplished in inappropriate regions rather than in the  
22 conserved ones. It must be outlined that though a certain sequence has a very low  
23 frequency, e. g. one, it may have been the dominant variant from an infected  
24 patient. This means that the unknown inpatient variability may modulate

1 behavior of virus resistant quasispecies turning the least frequent variant into the  
2 dominant one. The fact is that conservation depended on the number of DNA  
3 codons that could result in a particular amino acid sequence, that is why in the  
4 YMMD motif the conservation was less than the expected. Curiously, regions 1 and  
5 7 participate in DNA and dNTP binding, with only 2 and 4 non-consecutive amino  
6 acids respectively. These a.a. are involved in the two tasks mentioned, and had the  
7 best scores with the least number of combinations, under 48, from which two of  
8 them occurred in more than 900 sequences. Together, using the group of  
9 combinations of two or more regions will help fighting against viral variants. These  
10 results were obtained with the alignment of whole group M plus recombinants,  
11 however, nucleotide variation observed in subtypes different from B is important  
12 because this variability is acceptable in terms of enzyme function. This means that  
13 irrespective of the subtype the mutation can happen, and under selective drug  
14 pressure has a higher probability of appearance, without meaning that this variant  
15 will convert to another subtype. Still, there is a difference between naive viral  
16 variants and resistant variants, which not only has to be addressed *in vitro*, but  
17 which was the reason we analyzed sequences from resistant variants. Since  
18 resistant variants under certain drug inhibitor share resistance mutations, without  
19 changing drug selective pressure these mutations will become selected and thus  
20 conserved in the resistant viral population. Based on this fact, we searched within  
21 the previously selected conserved regions, resistant mutations that could increase  
22 the number of conserved consecutive amino acids in order to test it for shRNAs  
23 guide strand design. In fact we found much more variability, for example for  
24 regimen ZDV-3TC-EFV which was composed of 1381 sequences, the most

1 frequent combination occurred 825 times. This is less than in the curated alignment  
2 of whole group M plus recombinants, which is much more divergent. These results  
3 evidence the variability capacity of HIV virus, which undoubtedly is an obstacle for  
4 HIV gene therapy in an infected patient. Thus, depending on the regimen, the  
5 conserved nucleotide position in multiple alignments for YMDD ranged from 1 to 8  
6 of the 12 positions (other regimen specific alignments were used, data not shown).  
7 Unfortunately, this range is in certain degree a matter of the number of sequences  
8 used in the alignment, which means we have not arrived to the minimum number of  
9 sequences that represent HIV variability. For the region with the least number of  
10 unique combinations, amino acid reference sequence HXB2 is WPLTEEK, which  
11 can be formed by 512 different nucleotide sequences. It could be less if there is a  
12 preference for codon usage in HIV virus. Looking at the reference alignment for pol  
13 polyprotein (HIV Databases) the following mutations were seen: W24R, P25LTS,  
14 T26SA, E27KAGR, E28K, K29ER, which results in 286,654,464 possibilities. Since  
15 a region of 21 nt corresponds to 7 amino acids, the reason why the region of the  
16 YMMD motif being the most conserved did not show as expected, the least number  
17 of unique combinations, lies in the nearby amino acids, which can be coded by  
18 several DNA codons. Other conserved regions with nearby amino acids codified by  
19 less number of DNA codons could render better results.

20       Recently, McIntyre et al. 2009 analyzed a big set of sequences and find 96  
21 shRNAs for maximal coverage of the virus genome [56]. Even if our work had a  
22 similar approach, the final aim is different since they chose one shRNA for region  
23 and not the most frequent variants for each of them, which will reduce variability,  
24 as we did. For each of the regions they found, there are surely other frequent

1 variants that should be taken into account. If not, having more regions will result in  
2 increased variability reducing the therapy effectiveness, by means of the selection  
3 of different naturally occurring variants plus the resistant ones. On the other hand,  
4 targeting accessory genes could help in reducing pathogenicity of the clinical  
5 infection, but not necessarily viral replication, even though some accessory genes  
6 have a role in essential steps of viral cycle [18, 20, 21, 120-122]. Since nef protein  
7 is important in the context of the infection, as well as other accessory genes, this  
8 could explain the efficiency of silencing HIV in a humanized mice model, perhaps  
9 they are needed for viral fitness [47].

10 Our target was specifically RT because it cannot be deleted from the  
11 genome, besides it is the main source of virus variability. Furthermore, it is one of  
12 the targets of HAART, and there is a lot of information about mutations that occur  
13 under selective drug pressure [19, 62, 69, 123, 124]. Reverse transcription is  
14 absolutely necessary for viral integration and infection of other cells, so its silencing  
15 would reduce not only variability, but also the number of cells that become actively  
16 and latently infected, and thus HIV cell reservoirs [18, 125]. We did not target other  
17 essential enzymes [122], because targeting various enzymes or regions  
18 simultaneously would add variability (as was explained before), increasing the  
19 number of shRNAs needed for an efficient silencing. Nonetheless, multiple  
20 conserved regions could be the better approach as it has been previously  
21 proposed [43, 51, 56, 119], but only if the most frequent combinations for each  
22 conserved region are used for a shRNAs combinatorial design, in order to be able  
23 to target the most frequent variants. They will not necessarily be targeted by all the

1 groups of shRNAs simultaneously, but targeting a viral variant by at least one  
2 shRNA could result in silencing.

3         It might be helpful to study sequences from different specific regimens in  
4 order to find the better subset of shRNAs that should silence the emergent  
5 resistant isolates. This would be preferable to changing the previous regimen after  
6 a HAART failure, which will only change the viral population that is being selected,  
7 increasing the probability of failure with each regimen change. On the other hand, it  
8 might be useful to evaluate other regions with our approach in order to find better  
9 Score values with less number of shRNAs molecules, and perhaps find along the  
10 viral genome a region with a set of sequences that can cover the greatest  
11 percentage of virus variants, like for example LTR, which was proved to be very  
12 conserved [56]. Full-length genome sequencing of resistant viruses would be  
13 convenient for detecting target regions either in *pol* or in other genes, which can  
14 also target sequences from therapy-naive patients.

15         Efficient silencing was achieved *in vivo* in a mouse model using shRNA  
16 targeting a region within *nef* gene [47] even when before it was seen *in vitro* that  
17 *nef* gene was deleted as a route of escape [116]. In the pre-made curated multiple  
18 alignments we found no conservation in any of the 21-nucleotide positions used for  
19 the *nef* sequence (data not shown). In other words, silencing was directed to  
20 certain variants of the virus. In humans, quasispecies may be an obstacle for  
21 silencing, and using shRNAs against few variants of a region, as an initial therapy,  
22 will select the untargeted variants or induce emergence of resistance. The same is  
23 true in highly “conserved regions” of HIV-1 that surprisingly, reveal unexpected  
24 variability, which was even more evident in resistant isolates, and structure

1 alignments. Since mainly all the structures used were in complex with a specific  
2 ligand, different to all of them, observed variability may be caused by flexibility of  
3 the enzyme or structural conformational motions in order to achieve binding. In  
4 addition, variability is also a result of individual polymorphisms' that do not count as  
5 resistant mutations, but help in conserving the whole structure. Curiously, the  
6 conserved regions under study showed less structural alignment than we  
7 expected. Probably, even without mutations within this study region,  
8 polymorphisms distant from important residues create smooth structural  
9 conformational changes that are advantageous and maintain functional activity  
10 undamaged, as has being proposed before [126, 127]. So based on amino acid  
11 variability of conserved regions observed in structural alignments, we expected  
12 nucleotide variability to be high as we found.

13 Sequences identified here, might result in very efficient silencing once  
14 incorporated to an efficient hairpin [128], depending on intrapatient variability.  
15 Difference in silencing will only be evident proving the sequences *in vitro* and *in*  
16 *vivo*. If the number of sequences introduced in the expression vector does not  
17 become a problem, sequences from region 2 could work well since they have the  
18 higher scores. On the other hand, having a less number of molecules targeting  
19 more number of variants could be better in terms of expression [43]. The fact is  
20 that the combination of sequences chosen for the shRNA design is crucial in HIV-1  
21 gene therapy for a life-long efficacy. We do not discard the possibility of finding  
22 even better sequences for shRNA design. For this, the script must be improved,  
23 and other parameters should be included.

1 Most Blast hits resulted in non-significant E-values probably because short  
2 sequences have higher probability of occurring in the database. Nonetheless, even  
3 if the E-values we obtained are generally not considered as significant it is  
4 necessary to be careful since the targets exist in the databases, and in most cases  
5 identity was 100%. It was favorable that most of the probable sequences or unique  
6 combinations did not had a 100% overlapping with blast hits, only 9 from ZDV-  
7 3TC-NVP had full 21 nt alignment but not identity. In this case overlapping is  
8 important since the size of 21 nts is necessary for the intracellular-silencing  
9 pathway, and because normally shRNA are designed to target a fully  
10 complementary sequence [129, 130]. This year Liu et al., proposed a partially  
11 complementary approach for HIV inhibition, using shRNAs and artificial miRNAs, in  
12 which targeting not perfect sites by these two types of molecules resulted in  
13 silencing: miRNA and shRNA were capable of inhibiting multiple imperfect targets  
14 in the 3' UTR (untranslated region) of the viral genome. Our ESTs blast results  
15 showed that the overlapping hit aligned usually in the center, in some cases  
16 completely at the 3' end or at the 5' end of the guide strand. In these cases, identity  
17 was 100%. Hit sequences were usually shorter, but they ranged from 16 to 21 nts,  
18 which could represent high-risk off-target effects over the human mRNA  
19 sequences, particularly if they are expressed by the hematopoietic cell line. On the  
20 other hand, mutations in viral sequences were seen in different positions with no  
21 apparent pattern. All the possible targets of a partially complementarity strategy  
22 should be analyzed against ESTs databases in order to be sure to reduce off-  
23 target effects, since partially complementary sequences open the possibility of  
24 targeting cell genes that other wise would not be targeted if complete



1 complementarity was required for silencing. In addition, the blast results, obtained  
2 with the blast search against Human EST databases should be also addressed *in*  
3 *vitro* in order to confirm that there is no silencing of these human mRNAs species,  
4 even though overlapping is fewer than 95%. Each entry should be carefully  
5 analyzed in order to identify the specific cells that express the mRNA, when in  
6 development is it expressed and if it is already silenced by natural miRNAs. Many  
7 of the blast hits were ESTs associated with different types of cancer, perhaps in  
8 these cases off-target effects could be advantageous (e.g. Burritt's Lymphoma).

9       Also, using a multiple alignment of only a specific subtype, e. p. subtype B,  
10 could bring better results, since it should be more conserved than the reference  
11 multiple alignment used in this study that includes sequences from whole group M  
12 including recombinants. Our results are encouraging towards a more exhaustive  
13 study of HIV variability, and towards the finding of much better shRNAs based on  
14 subtype specific multiple alignments, which could be useful for a region based  
15 attempt. The fact that RT-Rtv domain is evolutionarily conserved among kingdoms,  
16 explains why it is relatively variable in HIV-1. Based on this argument, the next step  
17 would be to target unique viral sequences like *tat* or *rev* genes in combination with  
18 conserved regions of *env* (gp 120), targeting most frequent variants, in order to  
19 reduce viral variability for immune system, and help it control the infection. A  
20 possible region could be the residues important for CXCR5 binding, which are  
21 hidden until CD4 receptor binding with gp 120 [131].

22       In another hypothetical approach shRNAs could be used as a complement  
23 to HAART, where shRNAs should be directed against the specific resistant variants  
24 that emerge during a specified regimen. In this case, while HAART inhibits the

1 dominant and highly fit variants, shRNAs could inhibit a high percentage of the less  
2 fit emergent resistant variants. Perhaps specific mutations could be induced  
3 voluntarily in order to select variants that can be easily targeted by gene therapy.  
4 Our approach resulted in template guide strand sequences that can be used for  
5 shRNA design in order to be tested alone or under the specified regimen selective  
6 pressure. In a future prospect, unless we can find a better small number of  
7 combinations of sequences that could target all the variants, the ideal situation  
8 should be to initiate both therapies. shRNAs should only be started soon after  
9 virologic failure is evident, since most mutations occur under drug selective  
10 pressure. Indeed, further studies are needed in order to find regimens whose  
11 resistance pattern can be specified and better controlled using the least number of  
12 shRNAs. Proving real *in vivo* efficiency in a combination therapy as well as  
13 identifying off-target and non-specific off-target effects of the shRNAs molecules is  
14 absolutely required before starting any approach.

15

16

## 17 **ACKNOWLEDGEMENTS**

18 This study was sponsored by Instituto Nacional de Salud (INS) in Bogotá D.C.,  
19 Colombia, and Departamento Administrativo de Ciencia Tecnología e Innovación  
20 COLCIENCIAS Bogotá D.C., Colombia. Authors thank Dr. Andrés Paez Martinez  
21 at Instituto Nacional de Salud in Bogotá, for his competent assistance in the  
22 preparation of this manuscript.

23

1 **DECLARATION OF CONFLICT OF INTERESTS**

2 None

3

4 **FINANCIAL SUPPORT**

5 Departamento Administrativo de Ciencia Tecnología e Innovación COLCIENCIAS,

6 Bogota D.C., Colombia.

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

1

2 **SUPPLEMENTARY MATERIAL**

**Table 1. Clinical data for genotyped samples**

---

No. Samples	Recluse	55
	Analyzed	25
	Genotyped	25
Gender	Males	20
	Females	5
Age	Range	23 - 65
	Mean	30
Viral Load copies/ml	Range	1040 - >750000
	Mean	290,890
Final Condition	Alive	20
	Dead	5

---

3

4

5

6

7

Table 2. Antiretroviral Resistance data

Class	Inhibitor	Resistance (number of sequences)			
		Potential low level	Low level	Intermediate	High level
NRTI	3TC	1	0	1	4
	ABC	3	0	2	1
	AZT	0	0	1	1
	D4T	1	1	1	1
	DDI	1	0	3	
	FTC	1	0	1	4
	TDF	0	0	3	
NNRTI	DLV	0	0	1	3
	EFV	0	0	0	4
	ETR	1	2	1	0
	NVP	0	0	0	4
PRI	ATV/r	3	0	0	0
	DRV/r	0	0	0	0
	FPV/r	3	0	0	0
	IDV/r	3	0	0	0
	LPV/r	3	0	0	0
	NFV	0	2	1	0
	SQV/r	0	2	0	0

TPV/r                    2                    0                    0                    0

---

3TC lamivudine, ABC abacabir, AZT zidovudine, D4T stavudine, DDI didanosine, FTC emtricitabine, TDF tenofovir, DLV delavirdine, EFV efavirenz, ETR etravirine, NVP nevirapine, ATV/r atazanavir, DRV/r darunavir, FPV/r fosamprenavir, IDV/r indinavir, LPV/r lopinavir, NFV nelfinavir, SQV saquinavir, TPV tripanavir.

Antiretrovirals ARV/r, ritonavir boosted.

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18

1   **REFERENCES**

2

3   1.    Cohen, J (2000). Searching for the epidemic's origins. *Science* **288**: 2164-2165.

4

5   1.    Cohen, J (2000). Searching for the epidemic's origins. *Science* **288**: 2164-2165.

6   2.    Cohen, MS (2006). Amplified transmission of HIV-1: missing link in the HIV  
7    pandemic. *Trans Am Clin Climatol Assoc* **117**: 213-224; discussion 225.

8   3.    Mazumder, A, Engelman, A, Craigie, R, Fesen, M, and Pommier, Y (1994).  
9    Intermolecular disintegration and intramolecular strand transfer activities of wild-  
10   type and mutant HIV-1 integrase. *Nucleic Acids Res* **22**: 1037-1043.

11  4.    Wu, Y (2004). HIV-1 gene expression: lessons from provirus and non-integrated  
12   DNA. *Retrovirology* **1**: 13

13  5.    Chiu, TK, and Davies, DR (2004). Structure and function of HIV-1 integrase. *Curr*  
14   *Top Med Chem* **4**: 965-977.

15  6.    Garber, DA, Silvestri, G, and Feinberg, MB (2004). Prospects for an AIDS vaccine:  
16   three big questions, no easy answers. *Lancet Infect Dis* **4**: 397-413.

17  7.    Tripathi, P, and Agrawal, S (2007). Immunobiology of human immunodeficiency  
18   virus infection. *Indian J Med Microbiol* **25**: 311-322.

19  8.    Potter, SJ, Chew, CB, Steain, M, Dwyer, DE, and Saksena, NK (2004). Obstacles to  
20   successful antiretroviral treatment of HIV-1 infection: problems & perspectives.  
21   *Indian J Med Res* **119**: 217-237.

22  9.    Sethi, AK (2004). Adherence and HIV drug resistance. *HIV Clin Trials* **5**: 112-115.

23  10.   Bennett, DE, Bertagnolio, S, Sutherland, D, and Gilks, CF (2008). The World  
24   Health Organization's global strategy for prevention and assessment of HIV drug  
25   resistance. *Antivir Ther* **13 Suppl 2**: 1-13.

- 1 11. Boasso, A, and Shearer, GM (2008). Chronic innate immune activation as a cause  
2 of HIV-1 immunopathogenesis. *Clin Immunol* **126**: 235-242.
- 3 12. Moriuchi, M, and Moriuchi, H (2002). In vitro reactivation of HIV-1 by stimulation  
4 with herpes simplex virus infection. *Sex Transm Dis* **29**: 308-309.
- 5 13. Weinberg, MS, and Morris, KV (2006). Are viral-encoded microRNAs mediating  
6 latent HIV-1 infection? *DNA Cell Biol* **25**: 223-231.
- 7 14. King, SR, Duggal, NK, Ndongmo, CB, Pacut, C, and Telesnitsky, A (2008).  
8 Pseudodiploid genome organization AIDS full-length human immunodeficiency  
9 virus type 1 DNA synthesis. *J Virol* **82**: 2376-2384.
- 10 15. Cornelissen, M, van den Burg, R, Zorgdrager, F, Lukashov, V, and Goudsmit, J  
11 (1997). pol gene diversity of five human immunodeficiency virus type 1 subtypes:  
12 evidence for naturally occurring mutations that contribute to drug resistance, limited  
13 recombination patterns, and common ancestry for subtypes B and D. *J Virol* **71**:  
14 6348-6358.
- 15 16. van der Kuyl, AC, and Cornelissen, M (2007). Identifying HIV-1 dual infections.  
16 *Retrovirology* **4**: 67.
- 17 17. Greatorex, J (2004). The retroviral RNA dimer linkage: different structures may  
18 reflect different roles. *Retrovirology* **1**: 22.
- 19 18. Freed, EO (2001). HIV-1 replication. *Somat Cell Mol Genet* **26**: 13-33.
- 20 19. Basu, VP, Song, M, Gao, L, Rigby, ST, Hanson, MN, and Bambara, RA (2008).  
21 Strand transfer events during HIV-1 reverse transcription. *Virus Res* **134**: 19-38.
- 22 20. Basavapathruni, A, and Anderson, KS (2007). Reverse transcription of the HIV-1  
23 pandemic. *FASEB J* **21**: 3795-3808.



- 1 21. Cochrane, AW, McNally, MT, and Mouland, AJ (2006). The retrovirus RNA  
2 trafficking granule: from birth to maturity. *Retrovirology* **3**: 18.
- 3 22. Martinez-Picado, J, and Martinez, MA (2008). HIV-1 reverse transcriptase inhibitor  
4 resistance mutations and fitness: a view from the clinic and ex vivo. *Virus Res* **134**:  
5 104-123.
- 6 23. Alimonti, JB, Ball, TB, and Fowke, KR (2003). Mechanisms of CD4+ T  
7 lymphocyte cell death in human immunodeficiency virus infection and AIDS. *J Gen  
8 Virol* **84**: 1649-1661.
- 9 24. Collins, KR, *et al.* (2002). Human immunodeficiency virus type 1 (HIV-1)  
10 quasispecies at the sites of Mycobacterium tuberculosis infection contribute to  
11 systemic HIV-1 heterogeneity. *J Virol* **76**: 1697-1706.
- 12 25. Quan, Y, Liang, C, Brenner, BG, and Wainberg, MA (2009). Multidrug-Resistant  
13 Variants of HIV Type 1 (HIV-1) Can Exist in Cells as Defective Quasispecies and  
14 Be Rescued by Superinfection with Other Defective HIV-1 Variants. *J Infect Dis*  
15 **200**: 1479-1483.
- 16 26. Quinones-Mateu, ME, and Arts, EJ (2002). Fitness of drug resistant HIV-1:  
17 methodology and clinical implications. *Drug Resist Updat* **5**: 224-233.
- 18 27. Nigam, PK, and Kerketta, M (2006). AIDS vaccine: present status and future  
19 challenges. *Indian J Dermatol Venereol Leprol* **72**: 8-18.
- 20 28. Heeney, JL, Dalgleish, AG, and Weiss, RA (2006). Origins of HIV and the  
21 evolution of resistance to AIDS. *Science* **313**: 462-466.
- 22 29. Holmes, EC (2001). On the origin and evolution of the human immunodeficiency  
23 virus (HIV). *Biol Rev Camb Philos Soc* **76**: 239-254.

- 1 30. Piacenti, FJ (2006). An update and review of antiretroviral therapy.  
2 *Pharmacotherapy* **26**: 1111-1133.
- 3 31. Cavalcanti, AM, Lacerda, HR, Brito, AM, Pereira, S, Medeiros, D, and Oliveira, S  
4 (2007). Antiretroviral resistance in individuals presenting therapeutic failure and  
5 subtypes of the human immunodeficiency virus type 1 in the Northeast Region of  
6 Brazil. *Mem Inst Oswaldo Cruz* **102**: 785-792.
- 7 32. Kanzaki, LI, Ornelas, SS, and Arganaraz, ER (2008). RNA interference and HIV-1  
8 infection. *Rev Med Virol* **18**: 5-18.
- 9 33. Marathe, JG, and Wooley, DP (2007). Is gene therapy a good therapeutic approach  
10 for HIV-positive patients? *Genet Vaccines Ther* **5**: 5.
- 11 34. Agrawal, N, Dasaradhi, PV, Mohammed, A, Malhotra, P, Bhatnagar, RK, and  
12 Mukherjee, SK (2003). RNA interference: biology, mechanism, and applications.  
13 *Microbiol Mol Biol Rev* **67**: 657-685.
- 14 35. Shafer, RW, Rhee, SY, and Bennett, DE (2008). Consensus drug resistance  
15 mutations for epidemiological surveillance: basic principles and potential  
16 controversies. *Antivir Ther* **13 Suppl 2**: 59-68.
- 17 36. Lucas, GM (2005). Treatment of experienced patients and antiretroviral resistance  
18 issues from the 12th CROI. *Hopkins HIV Rep* **17**: 1-3, 12.
- 19 37. Bhavan, KP, Kampalath, VN, and Overton, ET (2008). The aging of the HIV  
20 epidemic. *Curr HIV/AIDS Rep* **5**: 150-158.
- 21 38. Leung, RK, and Whittaker, PA (2005). RNA interference: from gene silencing to  
22 gene-specific therapeutics. *Pharmacol Ther* **107**: 222-239.
- 23 39. Tijsterman, M, Ketting, RF, and Plasterk, RH (2002). The genetics of RNA  
24 silencing. *Annu Rev Genet* **36**: 489-519.

- 1 40. Saayman, S, Barichievy, S, Capovilla, A, Morris, KV, Arbuthnot, P, and Weinberg,  
2 MS (2008). The efficacy of generating three independent anti-HIV-1 siRNAs from a  
3 single U6 RNA Pol III-expressed long hairpin RNA. *PLoS One* **3**: e2602.
- 4 41. Scherer, L, Rossi, JJ, and Weinberg, MS (2007). Progress and prospects: RNA-  
5 based therapies for treatment of HIV infection. *Gene Ther* **14**: 1057-1064.
- 6 42. Sano, M, Li, H, Nakanishi, M, and Rossi, JJ (2008). Expression of long anti-HIV-1  
7 hairpin RNAs for the generation of multiple siRNAs: advantages and limitations.  
8 *Mol Ther* **16**: 170-177.
- 9 43. Liu, YP, *et al.* (2009). Combinatorial RNAi Against HIV-1 Using Extended Short  
10 Hairpin RNAs. *Mol Ther*.
- 11 44. Miyano-Kurosaki, N, and Takaku, H (2006). Gene silencing of virus replication by  
12 RNA interference. *Handb Exp Pharmacol*: 151-171.
- 13 45. ter Brake, O, Konstantinova, P, Ceylan, M, and Berkhout, B (2006). Silencing of  
14 HIV-1 with RNA interference: a multiple shRNA approach. *Mol Ther* **14**: 883-892.
- 15 46. Kirchhoff, F (2008). Silencing HIV-1 In Vivo. *Cell* **134**: 566-568.
- 16 47. ter Brake, O, *et al.* (2009). Evaluation of safety and efficacy of RNAi against HIV-1  
17 in the human immune system (Rag-2(-/-)gammac(-/-)) mouse model. *Gene Ther* **16**:  
18 148-153.
- 19 48. Touraine, JL, Sanhadji, K, and Sembeil, R (2003). Gene therapy for human  
20 immunodeficiency virus infection in the humanized SCID mouse model. *Isr Med*  
21 *Assoc J* **5**: 863-867.
- 22 49. Shacklett, BL (2008). Can the new humanized mouse model give HIV research a  
23 boost. *PLoS Med* **5**: e13.

- 1 50. von Eije, KJ, and Berkhout, B (2009). RNA-interference-based gene therapy  
2 approaches to HIV type-1 treatment: tackling the hurdles from bench to bedside.  
3 *Antivir Chem Chemother* **19**: 221-233.
- 4 51. von Eije, KJ, ter Brake, O, and Berkhout, B (2008). Human immunodeficiency  
5 virus type 1 escape is restricted when conserved genome sequences are targeted by  
6 RNA interference. *J Virol* **82**: 2895-2903.
- 7 52. Tsygankov, AY (2009). Current developments in anti-HIV/AIDS gene therapy.  
8 *Curr Opin Investig Drugs* **10**: 137-149.
- 9 53. ter Brake, O, t Hooft, K, Liu, YP, Centlivre, M, von Eije, KJ, and Berkhout, B  
10 (2008). Lentiviral vector design for multiple shRNA expression and durable HIV-1  
11 inhibition. *Mol Ther* **16**: 557-564.
- 12 54. von Eije, KJ, ter Brake, O, and Berkhout, B (2009). Stringent testing identifies  
13 highly potent and escape-proof anti-HIV short hairpin RNAs. *J Gene Med* **11**: 459-  
14 467.
- 15 55. Naito, Y, *et al.* (2007). Optimal design and validation of antiviral siRNA for  
16 targeting HIV-1. *Retrovirology* **4**: 80.
- 17 56. McIntyre, GJ, Groneman, JL, Yu, YH, Jaramillo, A, Shen, S, and Applegate, TL  
18 (2009). 96 shRNAs designed for maximal coverage of HIV-1 variants.  
19 *Retrovirology* **6**: 55.
- 20 57. Joseph, A, Sango, K, and Goldstein, H (2009). Novel mouse models for  
21 understanding HIV-1 pathogenesis. *Methods Mol Biol* **485**: 311-327.
- 22 58. Koyanagi, Y, Tanaka, Y, Ito, M, and Yamamoto, N (2008). Humanized mice for  
23 human retrovirus infection. *Curr Top Microbiol Immunol* **324**: 133-148.

- 1 59. Zhang, L, Kovalev, GI, and Su, L (2007). HIV-1 infection and pathogenesis in a  
2 novel humanized mouse model. *Blood* **109**: 2978-2981.
- 3 60. Rodriguez-Arenas, MA, *et al.* (2006). Delay in the initiation of HAART, poorer  
4 virological response, and higher mortality among HIV-infected injecting drug users  
5 in Spain. *AIDS Res Hum Retroviruses* **22**: 715-723.
- 6 61. Richman, DD, *et al.* (2004). The prevalence of antiretroviral drug resistance in the  
7 United States. *AIDS* **18**: 1393-1401.
- 8 62. Baldwin, C, and Berkhout, B (2007). HIV-1 drug-resistance and drug-dependence.  
9 *Retrovirology* **4**: 78.
- 10 63. Berkhout, B (2009). Toward a durable anti-HIV gene therapy based on RNA  
11 interference. *Ann N Y Acad Sci* **1175**: 3-14.
- 12 64. Shafer, RW (2006). Rationale and uses of a public HIV drug-resistance database. *J*  
13 *Infect Dis* **194 Suppl 1**: S51-58.
- 14 65. Rhee, SY, Gonzales, MJ, Kantor, R, Betts, BJ, Ravela, J, and Shafer, RW (2003).  
15 Human immunodeficiency virus reverse transcriptase and protease sequence  
16 database. *Nucleic Acids Res* **31**: 298-303.
- 17 66. Shafer, RW, Jung, DR, Betts, BJ, Xi, Y, and Gonzales, MJ (2000). Human  
18 immunodeficiency virus reverse transcriptase and protease sequence database.  
19 *Nucleic Acids Res* **28**: 346-348.
- 20 67. Ding, J, Hughes, SH, and Arnold, E (1997). Protein-nucleic acid interactions and  
21 DNA conformation in a complex of human immunodeficiency virus type 1 reverse  
22 transcriptase with a double-stranded DNA template-primer. *Biopolymers* **44**: 125-  
23 138.

- 1 68. Marchler-Bauer, A, *et al.* (2007). CDD: a conserved domain database for interactive  
2 domain family analysis. *Nucleic Acids Res* **35**: D237-240.
- 3 69. Back, NK, *et al.* (1996). Reduced replication of 3TC-resistant HIV-1 variants in  
4 primary cells due to a processivity defect of the reverse transcriptase enzyme.  
5 *EMBO J* **15**: 4040-4049.
- 6 70. Coghill, P, Finn, RD, and Bateman, A (2008). Identifying protein domains with the  
7 Pfam database. *Curr Protoc Bioinformatics* **Chapter 2**: Unit 2 5.
- 8 71. Finn, R, Griffiths-Jones, S, and Bateman, A (2003). Identifying protein domains  
9 with the Pfam database. *Curr Protoc Bioinformatics* **Chapter 2**: Unit 2 5.
- 10 72. Finn, RD, *et al.* (2008). The Pfam protein families database. *Nucleic Acids Res* **36**:  
11 D281-288.
- 12 73. Hunter, S, *et al.* (2009). InterPro: the integrative protein signature database. *Nucleic*  
13 *Acids Res* **37**: D211-215.
- 14 74. Mulder, N, and Apweiler, R (2007). InterPro and InterProScan: tools for protein  
15 sequence classification and comparison. *Methods Mol Biol* **396**: 59-70.
- 16 75. Gifford, RJ, *et al.* (2009). The calibrated population resistance tool: standardized  
17 genotypic estimation of transmitted HIV-1 drug resistance. *Bioinformatics* **25**:  
18 1197-1198.
- 19 76. Rhee, SY, *et al.* (2006). HIV-1 pol mutation frequency by subtype and treatment  
20 experience: extension of the HIVseq program to seven non-B subtypes. *AIDS* **20**:  
21 643-651.
- 22 77. Berman, HM, *et al.* (2000). The Protein Data Bank. *Nucleic Acids Res* **28**: 235-242.
- 23 78. Berman, H, Henrick, K, and Nakamura, H (2003). Announcing the worldwide  
24 Protein Data Bank. *Nat Struct Biol* **10**: 980.

- 1 79. Das, K, *et al.* (2008). High-resolution structures of HIV-1 reverse  
2 transcriptase/TMC278 complexes: strategic flexibility explains potency against  
3 resistance mutations. *Proc Natl Acad Sci U S A* **105**: 1466-1471.
- 4 80. Zhao, Z, *et al.* (2008). Novel indole-3-sulfonamides as potent HIV non-nucleoside  
5 reverse transcriptase inhibitors (NNRTIs). *Bioorg Med Chem Lett* **18**: 554-559.
- 6 81. Ren, J, *et al.* (2008). Structural basis for the improved drug resistance profile of new  
7 generation benzophenone non-nucleoside HIV-1 reverse transcriptase inhibitors. *J*  
8 *Med Chem* **51**: 5000-5008.
- 9 82. Ren, J, *et al.* (2007). Relationship of potency and resilience to drug resistant  
10 mutations for GW420867X revealed by crystal structures of inhibitor complexes for  
11 wild-type, Leu100Ile, Lys101Glu, and Tyr188Cys mutant HIV-1 reverse  
12 transcriptases. *J Med Chem* **50**: 2301-2309.
- 13 83. Hopkins, AL, *et al.* (2004). Design of non-nucleoside inhibitors of HIV-1 reverse  
14 transcriptase with improved drug resistance properties. 1. *J Med Chem* **47**: 5912-  
15 5922.
- 16 84. Freeman, GA, *et al.* (2004). Design of non-nucleoside inhibitors of HIV-1 reverse  
17 transcriptase with improved drug resistance properties. 2. *J Med Chem* **47**: 5923-  
18 5936.
- 19 85. Ren, J, Nichols, CE, Chamberlain, PP, Weaver, KL, Short, SA, and Stammers, DK  
20 (2004). Crystal structures of HIV-1 reverse transcriptases mutated at codons 100,  
21 106 and 108 and mechanisms of resistance to non-nucleoside inhibitors. *J Mol Biol*  
22 **336**: 569-578.
- 23 86. Bauman, JD, *et al.* (2008). Crystal engineering of HIV-1 reverse transcriptase for  
24 structure-based drug design. *Nucleic Acids Res* **36**: 5083-5092.

- 1 87. Das, K, *et al.* (2004). Roles of conformational and positional adaptability in  
2 structure-based design of TMC125-R165335 (etravirine) and related non-nucleoside  
3 reverse transcriptase inhibitors that are highly potent and effective against wild-type  
4 and drug-resistant HIV-1 variants. *J Med Chem* **47**: 2550-2560.
- 5 88. Ren, J, Milton, J, Weaver, KL, Short, SA, Stuart, DI, and Stammers, DK (2000).  
6 Structural basis for the resilience of efavirenz (DMP-266) to drug resistance  
7 mutations in HIV-1 reverse transcriptase. *Structure* **8**: 1089-1094.
- 8 89. Das, K, Sarafianos, SG, Clark, AD, Jr., Boyer, PL, Hughes, SH, and Arnold, E  
9 (2007). Crystal structures of clinically relevant Lys103Asn/Tyr181Cys double  
10 mutant HIV-1 reverse transcriptase in complexes with ATP and non-nucleoside  
11 inhibitor HBY 097. *J Mol Biol* **365**: 77-89.
- 12 90. Shatsky, M, Nussinov, R, and Wolfson, HJ (2004). A method for simultaneous  
13 alignment of multiple protein structures. *Proteins* **56**: 143-156.
- 14 91. Guex, N, and Peitsch, MC (1997). SWISS-MODEL and the Swiss-PdbViewer: an  
15 environment for comparative protein modeling. *Electrophoresis* **18**: 2714-2723.
- 16 92. Schwede, T, Kopp, J, Guex, N, and Peitsch, MC (2003). SWISS-MODEL: An  
17 automated protein homology-modeling server. *Nucleic Acids Res* **31**: 3381-3385.
- 18 93. Shafer, RW, Stevenson, D, and Chan, B (1999). Human Immunodeficiency Virus  
19 Reverse Transcriptase and Protease Sequence Database. *Nucleic Acids Res* **27**: 348-  
20 352.
- 21 94. Kantor, R, Machekano, R, Gonzales, MJ, Dupnik, K, Schapiro, JM, and Shafer, RW  
22 (2001). Human Immunodeficiency Virus Reverse Transcriptase and Protease  
23 Sequence Database: an expanded data model integrating natural language text and  
24 sequence analysis programs. *Nucleic Acids Res* **29**: 296-299.



- 1 95. Katoh, K, Misawa, K, Kuma, K, and Miyata, T (2002). MAFFT: a novel method for  
2 rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids*  
3 *Res* **30**: 3059-3066.
- 4 96. Granados, CAD, Álvarez, C, Prada, G, Martínez, FL, and Sarmiento, CA (2006).  
5 Guía para el manejo de VIH/sida Basada en la evidencia In *Modelo de gestión*  
6 *programática en VIH/sida*. Programa de Apoyo a la Reforma de Salud Ajuste a los  
7 Planes de Beneficio y la Unidad de Pago por Capitación. MPS, COLOMBIA.
- 8 97. Liu, TF, and Shafer, RW (2006). Web resources for HIV type 1 genotypic-  
9 resistance test interpretation. *Clin Infect Dis* **42**: 1608-1618.
- 10 98. Reynolds, A, Leake, D, Boese, Q, Scaringe, S, Marshall, WS, and Khvorova, A  
11 (2004). Rational siRNA design for RNA interference. *Nat Biotechnol* **22**: 326-330.
- 12 99. Jagla, B, *et al.* (2005). Sequence characteristics of functional siRNAs. *RNA* **11**: 864-  
13 872.
- 14 100. Marchler-Bauer, A, *et al.* (2009). CDD: specific functional annotation with the  
15 Conserved Domain Database. *Nucleic Acids Res* **37**: D205-210.
- 16 101. Huang, H, Chopra, R, Verdine, GL, and Harrison, SC (1998). Structure of a  
17 covalently trapped catalytic complex of HIV-1 reverse transcriptase: implications  
18 for drug resistance. *Science* **282**: 1669-1675.
- 19 102. Prajapati, DG, Ramajayam, R, Yadav, MR, and Giridhar, R (2009). The search for  
20 potent, small molecule NNRTIs: A review. *Bioorg Med Chem* **17**: 5744-5762.
- 21 103. Eddy, SR (1998). Profile hidden Markov models. *Bioinformatics* **14**: 755-763.
- 22 104. Gong, W, *et al.* (2006). Integrated siRNA design based on surveying of features  
23 associated with high RNAi effectiveness. *BMC Bioinformatics* **7**: 516.

- 1 105. Ren, Y, Gong, W, Zhou, H, Wang, Y, Xiao, F, and Li, T (2009). siRecords: a  
2 database of mammalian RNAi experiments and efficacies. *Nucleic Acids Res* **37**:  
3 D146-149.
- 4 106. Tilesi, F, Fradiani, P, Socci, V, Willems, D, and Ascenzioni, F (2009). Design and  
5 validation of siRNAs and shRNAs. *Curr Opin Mol Ther* **11**: 156-164.
- 6 107. Bartel, DP (2004). MicroRNAs: genomics, biogenesis, mechanism, and function.  
7 *Cell* **116**: 281-297.
- 8 108. Zhang, C (2008). MicroRNomics: a newly emerging approach for disease biology.  
9 *Physiol Genomics* **33**: 139-147.
- 10 109. Zhang, B, Wang, Q, and Pan, X (2007). MicroRNAs and their regulatory roles in  
11 animals and plants. *J Cell Physiol* **210**: 279-289.
- 12 110. John, B, Enright, AJ, Aravin, A, Tuschl, T, Sander, C, and Marks, DS (2004).  
13 Human MicroRNA targets. *PLoS Biol* **2**: e363.
- 14 111. Miranda, KC, *et al.* (2006). A pattern-based method for the identification of  
15 MicroRNA binding sites and their corresponding heteroduplexes. *Cell* **126**: 1203-  
16 1217.
- 17 112. Liu, YP, Gruber, J, Haasnoot, J, Konstantinova, P, and Berkhout, B (2009). RNAi-  
18 mediated inhibition of HIV-1 by targeting partially complementary viral sequences.  
19 *Nucleic Acids Res.*
- 20 113. Boden, D, Pusch, O, and Ramratnam, B (2007). Overcoming HIV-1 resistance to  
21 RNA interference. *Front Biosci* **12**: 3104-3116.
- 22 114. Boden, D, Pusch, O, Lee, F, Tucker, L, and Ramratnam, B (2003). Human  
23 immunodeficiency virus type 1 escape from RNA interference. *J Virol* **77**: 11531-  
24 11535.

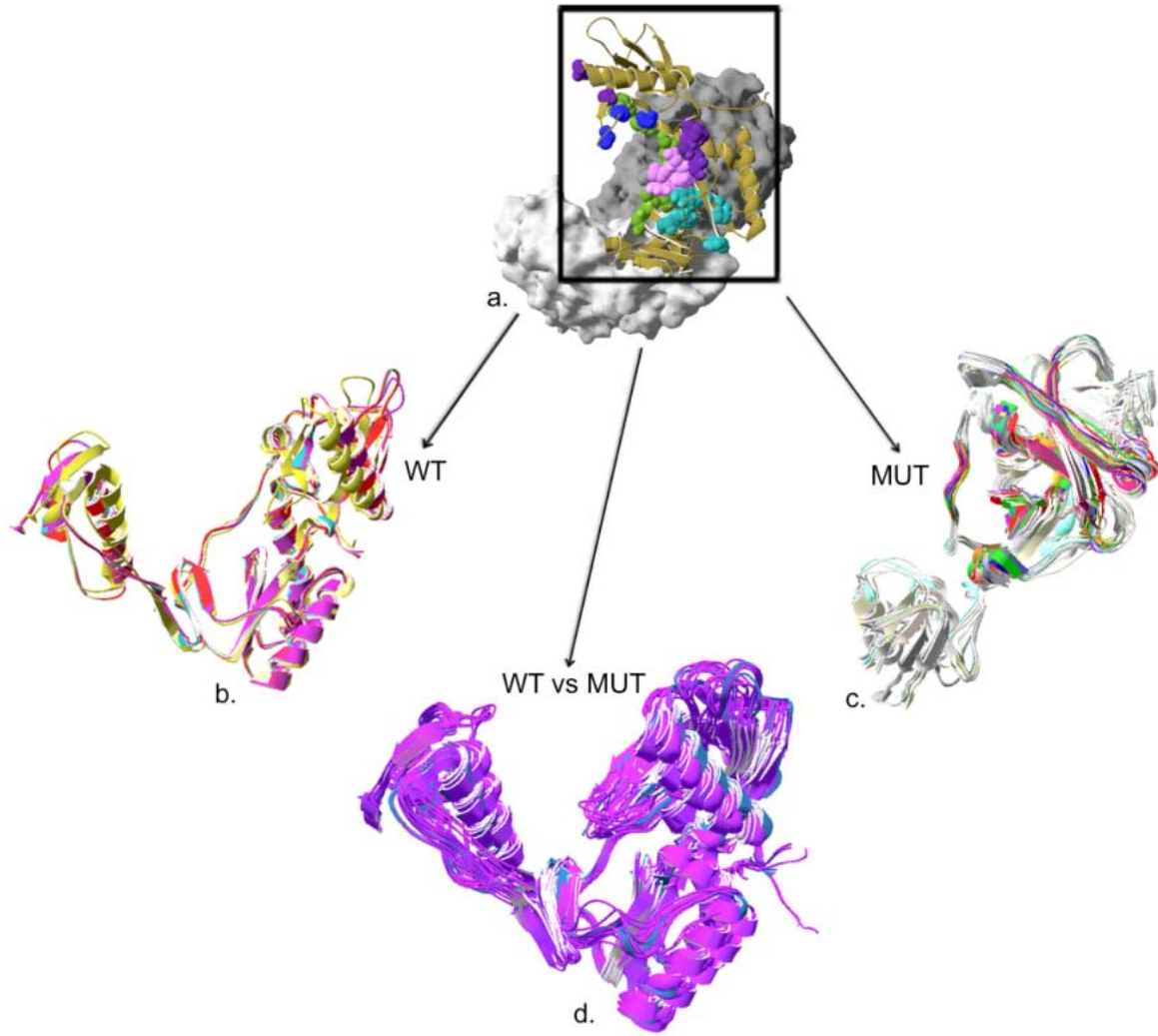
- 1 115. Boden, D, Pusch, O, and Ramratnam, B (2004). HIV-1-specific RNA interference.  
2 *Curr Opin Mol Ther* **6**: 373-380.
- 3 116. Das, AT, *et al.* (2004). Human immunodeficiency virus type 1 escapes from RNA  
4 interference-mediated inhibition. *J Virol* **78**: 2601-2605.
- 5 117. Westerhout, EM, Ooms, M, Vink, M, Das, AT, and Berkhout, B (2005). HIV-1 can  
6 escape from RNA interference by evolving an alternative structure in its RNA  
7 genome. *Nucleic Acids Res* **33**: 796-804.
- 8 118. Liu, YP, Haasnoot, J, and Berkhout, B (2007). Design of extended short hairpin  
9 RNAs for HIV-1 inhibition. *Nucleic Acids Res* **35**: 5683-5693.
- 10 119. Liu, YP, Haasnoot, J, ter Brake, O, Berkhout, B, and Konstantinova, P (2008).  
11 Inhibition of HIV-1 by multiple siRNAs expressed from a single microRNA  
12 polycistron. *Nucleic Acids Res* **36**: 2811-2824.
- 13 120. Bennasser, Y, and Jeang, KT (2006). HIV-1 Tat interaction with Dicer: requirement  
14 for RNA. *Retrovirology* **3**: 95.
- 15 121. Costin, JM (2007). Cytopathic mechanisms of HIV-1. *Virol J* **4**: 100.
- 16 122. Klimas, N, Koneru, AO, and Fletcher, MA (2008). Overview of HIV. *Psychosom*  
17 *Med* **70**: 523-530.
- 18 123. Bakhouch, K, *et al.* (2009). The prevalence of resistance-associated mutations to  
19 protease and reverse transcriptase inhibitors in treatment-naive (HIV1)-infected  
20 individuals in Casablanca, Morocco. *J Infect Dev Ctries* **3**: 380-391.
- 21 124. Boyle, BA, *et al.* (2008). Update on antiretroviral therapy: the 15th CROI. *AIDS*  
22 *Read* **18**: 273-278, C273.
- 23 125. Brenner, B, *et al.* (2004). Persistence of multidrug-resistant HIV-1 in primary  
24 infection leading to superinfection. *AIDS* **18**: 1653-1660.

- 1 126. Gutteridge, A, and Thornton, J (2004). Conformational change in substrate binding,  
2 catalysis and product release: an open and shut case? *FEBS Lett* **567**: 67-73.
- 3 127. Tousseignant, A, and Pelletier, JN (2004). Protein motions promote catalysis. *Chem*  
4 *Biol* **11**: 1037-1042.
- 5 128. Ding, H, Liao, G, Wang, H, and Zhou, Y (2007). Asymmetrically designed siRNAs  
6 and shRNAs enhance the strand specificity and efficacy in RNAi. *J RNAi Gene*  
7 *Silencing* **4**: 269-280.
- 8 129. Rao, DD, Vorhies, JS, Senzer, N, and Nemunaitis, J (2009). siRNA vs. shRNA:  
9 similarities and differences. *Adv Drug Deliv Rev* **61**: 746-759.
- 10 130. Cohen, MS, Hellmann, N, Levy, JA, DeCock, K, and Lange, J (2008). The spread,  
11 treatment, and prevention of HIV-1: evolution of a global pandemic. *J Clin Invest*  
12 **118**: 1244-1254.

13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24

1 **Figure 1.**

2



3

4

5

6

7

8

1 Fig1. a. RT crystallographic structure 2ZD1 (1.8 Å) highlighting inside the squared  
2 frame the residues within selected regions. Dark gray p66 subunit, in light gray  
3 p55, in dark blue active site residues involved dNTP binding (K65, R72, D110,  
4 V111, G112, D113, A114, Y115, Q151), in green active site residues involved in  
5 DNA Binding (L74, V75, D76, R78, N81, E89, Q91, L92, I94, G152, K154, P157,  
6 M230, G231), in purple active site residues with no specific annotations (W24, P25,  
7 F61), in pink YMDD motif (Y183, M184, D185, D186), and in light blue residues  
8 involved in NNRTI binding (L100, K101, K102, K103, V179, Y188, G190, F227; not  
9 conserved). Ribbon shows continuity between amino-acid chains. b. Active site  
10 structural alignment for 10 different RT wild type structures, in ribbon (RMS 0.81).  
11 Bound ligands and heteroatoms (hidden) were not taken into account for the  
12 structural alignment. Subunit p51 (A chain) and residues 297 to 545, are hidden,  
13 but aligned. Reference layer: structure 2ZD1. b. Structural alignment for 19  
14 different RT mutant enzymes, in ribbon. Selected regions are shown in highlighted  
15 colors, other not important residues in grayscale. Only 590 residues were  
16 structurally aligned. Structural alignment was generated with MultiProt [90].

**Table 1. Conserved residues and selected regions chosen for the analysis, their position in HXB2 reference genome, the important residues within them including their position with respect to RT enzyme (HXB2 numbering system), and the nucleotide sequence for each region.**

Region	Residue No.	Amino acid	Function	DNA	Evaluated region*
				codon in HXB2 genome	
1	24	W	DNA Binding	2619	<sup>2616</sup> caaTGGC AAttgaca <sub>2624</sub>
	25	P	and Active Site	2622	
2	61	F	Active Site	2730	<sup>2814</sup> TTTgccataaagAAA <sub>2831</sub>

		dNTP	
		Binding	
65	K	and	2742
		Active	
		Site	

---

		dNTP	
		Binding	
72	R	and	2763
		Active	
		Site	
<b>3</b>		DNA	
74	L		2769
75	V	Binding	2772
76	D	and	2775
78	R	Active	2781
81	N	Site	2790

<sup>2763</sup>AGAaaaTTAGTAGATttcAGAgaaacttAAT<sub>2792</sub>



	89	E	DNA	2814	
	91	Q	Binding	2820	
<b>4</b>	92	L	and	2823	<sup>2814</sup> GAAgttCAATTAggaATA <sub>2831</sub>
	94	I	Active Site	2829	
<hr/>					
	100	L		2847	
			NNRTI		
<b>5</b>	101	K	Binding	2850	
	102	K	site	2853	<sup>2847</sup> TTAAAAAAGAAA <sub>2859</sub>
	103	K		2856	
<hr/>					
	110	D		2877	
			dNTP		
	111	V	Binding	2880	
<b>6</b>	112	G	and	2883	
	113	D	Active	2886	<sup>2877</sup> GATGTGGGTGATGCATAT <sub>2894</sub>
	114	A	Site	2889	
	115	Y		2892	

			dNTP		
			Binding		
	151	Q	and	3000	
			Active		
<b>7</b>			Site		
	152	G	DNA	3003	<sup>3000</sup> CAGGGAtggAAAggatcaCCA <sub>3020</sub>
	154	K	Binding	3009	
			and		
	157	P	Active	3018	
			Site		

---

			NNRTI		
	179	V	Binding	3084	
<b>8</b>			site		<sup>3084</sup> GTTatctactatcaaTACAGTGATGATtTgTATgtaGGA <sub>3119</sub>
	183	Y	NNRTI	3096	
			Binding		

---

		Site, DNA Binding and Active Site	
184	M	variable dNTP Binding	3099
185	D	and Active Site	3102
186	D	Active Site	3105
188	Y	NNRTI	3111

			Binding site		
			NNRTI		
	190	G	Binding site	3117	
<hr/>					
			NNRTI		
	227	F	Binding site	3228	
	230	M	DNA	3237	
<b>9</b>			Binding and Active Site		
	231	G		3240	
					<sup>3228</sup> TTCcttggATGGGTtatgaactC <sub>3251</sub>

\*DNA codons in capslock correspond to the residues in the region. Nucleotide sequences correspond to the HXB2 reference genome, positions according to the HXB2 numbering system (map of coordinates). Each

nucleotide position will be considered the starting position of a window from which the following upstream 20  
nts are included

1 **Table 2.**

**Table 2. Prevalence of Resistant Mutations that reside inside selected regions**

Region <sup>1</sup>	HXB2		Mutations Prevalence		
	Position	Residue(s)	NRTI <sup>2</sup>	NNRTI <sup>3</sup>	RTI <sup>4</sup>
<b>2</b>	60	V	V (12)	V (16)	V (14)
	61	F	-	-	-
	62	A	V (2.2)	V (4.8)	V (14)
	63	I	-	-	-
	64	K	R (1.6)	R (2.1), H (1.2)	R (1.9)
	65	K	R (1.3)	R (2.6)	R (2.1)
	66	K	-	-	-
	67	D	N (26), G (1.1)	K (2.6), G (3.3)	N (38), G (2.5)
<b>4</b>	89	E	0		
	90	V	I (3.7)	I (2.4)	I (3.3)
	91-94	Q, L, G, I	-	-	-

<b>7</b>	151	Q	M(4.0)	M(2.3)	M(3.4)
	152-157	G, W, K, G, S, P	-	-	-
<b>8</b>	178	I	M(7.8), L(6.7)	M(7.4), L(6.2)	M(7.5), L(6.4)
	179	V	I(3.3)	I(9.6), D(1.7)	I(7.0), D(1.3)
	180	I	-	-	-
	181	Y	-	C(23), I(1.1)	C(14)
	182	Q, Y	-	-	-
	184	M	V(49)	V(50), I(1.9)	V(50), I(1.4)
	185	D, D, L	-	-	-
	188	Y	-	L(5.9)	L(3.5)
	189	V	-	I(2.3)	I(1.5)
	190	G	-	A(16), S(2.9)	
<b>9</b>	227	F	-	L(2.5)	L(1.6)
	228	L	-	H(17), R(5.9)	H(12), R(4.2)
	229	W	-	-	-

230	M <sup>1</sup>	-	-	L (1.8)
231	G <sup>2</sup>	-	-	-

---

<sup>1</sup>Corresponds to the numeration given to selected regions from table 1.

<sup>4</sup>NRTI prevalence was calculated from 3583 sequences exposed to these type of drugs

<sup>4</sup>NNRTI prevalence was calculated form 5070 sequences exposed to these type of drugs

<sup>4</sup>RTI prevalence was calculated from 17167 sequences exposed to either type of these drugs

<sup>4,3,4</sup>Mutations prevalence data is available at HIV Drug Resistance Database

<sup>1</sup>Residues directly involved in enzyme activity

In shade, selected regions.

---

1  
2



1  
2  
3  
4

**Table 3.**

**Table 3. Scores, unique combinations (guide strand template) and sequences for the regions selected for silencing**

Alignment	Region of study	Region in HXB2	Window	No. of U.C <sup>2</sup>	U.C <sup>2</sup>	Sequence	Score	Frequency	Score Conditions <sup>3</sup>																		
									1	2	3	4	5	6	7	8	9	10	11	12	13						
Curated Alignment Group M plus Recombinants	1	2616 - 2630	3234	46	1	AGGAGCAGATGATACAGTATT	8	959	N	Y	Y	Y	Y	N	Y	N	Y	Y	Y	Y	Y	Y	Y				
					2	AGGAGCAGATGATACAGTGTT	5	40	N	Y	N	N	Y	N	Y	N	Y	N	Y	Y	Y	Y	Y	Y	Y		
					3	AGGAGCAGATGACACAGTATT	8	32	N	Y	Y	Y	Y	N	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
					4	AGGGGCAGATGATACAGTATT	8	29	N	Y	Y	Y	Y	N	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
					5	AGGAGCAGATGATACAGTACT	7	24	N	Y	Y	Y	Y	N	Y	N	Y	Y	N	Y	Y	N	Y	Y	Y	Y	Y
		1	AGCAGATGATACAGTATTAGA	9	942	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	N	Y	Y	Y	Y			
		2	AGCAGATGATACAGTGTTAGA	8	40	N	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y		
		3	AGCAGATGACACAGTATTAGA	8	32	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y		
		4	GGCAGATGATACAGTATTAGA	9	27	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y		
		5	AGCAGATGATACAGTACTAGA	8	23	N	Y	Y	Y	Y	Y	Y	N	Y	Y	N	Y	Y	N	Y	Y	Y	Y	Y	Y		
	7	3000 - 3014	3648	47	1	CCAGTATTTGCCATAAAAAAG	10	167	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y			
	2				CCAATATTTGCCATAAAAAAG	10	161	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y		
	3				CCAGTATTTGCTATAAAAAAG	10	27	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
	4				CCAATATTTGCTATAAAAAAG	10	25	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
	5				CCAGTATTTGCCATAAAGAAA	9	255	Y	Y	N	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
	6				CCAGTATTTGCTATAAAGAAA	9	165	Y	Y	N	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
	7				CCAATATTTGCCATAAAGAAA	9	113	Y	Y	N	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
	8				CCAATATTTGCTATAAAGAAA	9	91	Y	Y	N	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
	9				CCAGTATTTGCCATAAAGAAA	9	25	Y	Y	N	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
	10				CCAATATTTGCAATAAAGAAA	9	22	Y	Y	N	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	

D4T-3TC-NVP (stavudine, lamivudine, nevirapine) 91 sequences	2	179	28	1	TATTTGCTATAAAGAAAAAGA	8	14	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	N	N	Y			
				2	TATTTGCTATAAAGAAAAAG	9	8	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	N	Y	
				3	TATTTGCCATAAAGAAAAAAG	9	20	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	N	Y
				4	TATTTGCTATAAAGAAAAAGG	8	16	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	N	N	Y
				5	TATTTGCCATAAAGAAAAAGG	8	7	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	N	N	Y
				6	TATTTGCCATAAAAAAGAAGG	7	6	N	Y	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	N	N	Y
	4	264	11	1	GGAAGTTCATTAGGAATACC	8	60	Y	N	Y	Y	Y	N	Y	Y	Y	Y	Y	N	Y	Y		
				2	GGAAGTTCAGTTAGGAATACC	8	7	Y	N	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	N	Y	Y	
				3	GGAAATTCATTAGGAATACC	7	7	Y	N	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	N	N	Y	
				4	GGAGGTTCATTAGGAATACC	8	6	Y	N	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	N	Y	Y	
	7	459	11	1	GAAAGGATCACCAGCAATATT	8	54	Y	N	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	N	Y		
				2	GAAAGGATCACCGCAATATT	7	14	Y	N	Y	Y	N	N	Y	Y	Y	Y	Y	Y	Y	N	Y	
				3	GAAGGATCACCAGCAATATT	9	8	Y	N	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	
	8	565	21	1	GTAGGATCTGACTTAGAAATA	9	35	Y	N	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y		
				2	GTAGGATCTGATTTAGAAATA	9	18	Y	N	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	
				3	GTAGGATCTGATTTAGAAATA	9	14	Y	N	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	
	ZDV-3TC-EFV (zidovudine, Lamivudine, efavirenz) 1381 sequences	1	72	74	1	GCCATTGACAGAAGAAAAAC	9	5	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y		
					2	GCCATTGACAGAAGAAAAAT	9	825	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
					3	GCCACTGACAGAAGAAAAAT	9	56	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
					4	GCCATTGACAGAAGAAAGAT	6	21	N	N	N	N	Y	Y	Y	Y	Y	N	Y	Y	Y	Y	Y
					5	GCCCTTGACAGAAGAAAAAT	9	20	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
					6	GCCATTGACAGAGAAAAAAT	8	42	N	N	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
					7	GCCATTAACAGAAGAGAAAAT	8	15	N	N	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y
					8	GCCATTGACGGAAAGAAAAAT	9	28	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
9					GCCATTGTCAGAAGAAAAAT	9	13	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
10					GCCATTGACAAAAGAAAAAT	8	11	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	N	Y	
11					GCCATTGACAGAAGAGAAAAT	8	144	N	N	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	
12					GCCATTAACAGAAGAAAAAT	9	33	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	
13					GCCATTGACAAAAGAGAAAAT	7	17	N	N	Y	Y	Y	Y	N	Y	Y	Y	Y	Y	Y	N	Y	
73	77	1	CCATTGACAGAAGAGAAAATT	8	31	Y	N	Y	Y	N	N	Y	Y	Y	Y	Y	Y	Y	Y				
		2	CCATTGTCAGAAGAAAAATA	9	13	Y	N	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y				

---

3	CCATTGACAGAAGAAAAGATA	8	21	Y	N	Y	Y	N	Y	Y	Y	N	Y	Y	Y	Y	Y	Y
4	CCATTGACAGAAGAGAAAATA	8	114	Y	N	Y	Y	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y
5	CCATTGACAAAAAGAAAATA	7	17	Y	N	Y	Y	N	N	Y	Y	Y	Y	Y	Y	N	Y	Y
6	CCATTGACAGAGGAAAAATA	9	42	Y	N	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
7	CCATTGACAGAAGAAAAAATT	9	13	Y	N	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
8	CCATTGACGGAAGAAAAATA	9	28	Y	N	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
9	CCCTTGACAGAAGAAAAATA	8	20	N	N	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
10	CCATTAACAGAAGAAAAATA	9	33	Y	N	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
11	CCATTGACAGAAGAAAAATA	9	813	Y	N	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
12	CCATTAACAGAAAGAAAATT	8	15	Y	N	Y	Y	N	N	Y	Y	Y	Y	Y	Y	Y	Y	Y
13	CCACTGACAGAAGAAAAATA	9	56	Y	N	Y	Y	N	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y

---

<sup>1</sup>M.A. Multiple Alignment

<sup>2</sup>U.C.= Unique Combinations (with respect the total No. of sequences in the M.A)

<sup>3</sup>Score conditions:

1. Has A in position 3
2. Has T in position 10
3. No C/G in position 19, but has A
4. No C/G in position 19
5. No G in position 13
6. Has A or T in position 15
7. Has A or T in position 16
8. Has A or T in position 17
9. Has A or T in position 18
10. Has A or T in position 19
11. Has A or T in position 20
12. No more than three consecutive A's from position 1 to 14
13. The sum of C or G is less than 11

Y = accomplished

N= not accomplished

1

2

1  
2

**TABLE 4. Most significant Blast Hits against UniSTS only hits in minus strands are shown**

Alignment	No.	HIT	U.C	E-value	Bit score	Overlap	Hit frame/strand
<b>Curated alignment</b>	1	NT_005403	R1,w2,c4				
	1	NT_005403	R1,w2,c5	0,3	36,18	76,19	Minus
	2	NW_921618					
	3	NT_016297	R2,w1,c3	0,3	36,18	76,19	Minus
	4	NW_922073					
	5	NT_025741	R2,w1,c6	<b>0,08</b>	38,16	80,95	Minus
	6	NW_923184					
	7	NT_009714					
	8	NT_016354	R2,w1,c8	0,3	36,18	76,19	Minus
	9	NW_925328					
	10	NW_922162					
	11	NT_025741	R2,w1,c9	0,3	36,18	76,19	Minus
	12	NW_923184					
	13	NT_022517					
	14	NT_004487	R2,w1,c10	0,3	36,18	76,19	Minus
	15	NW_921651					
16	NW_926128						

Blast search against Uni STS database was only perform with Curated Alignment combinations .

3  
4  
5  
6  
7  
8  
9

Table 5. Blast Most significant Hits against Human ESTs database

Alignment	Region	HIT	U.C	E-value	Bit score	Overlap
<b>Curated</b>	7	BF240717	R7, 8	0,21231	36,1753	85,71
		DA099283	R7, 4 R7, 8	0,838917	34,1929	
		AW503859	R7, 4 R7, 8	0,838917	34,1929	80,95
		D17012	R7, 4 R7, 8	0,838917	34,1929	
<b>D4T-3TC-NVP (stavudine- lamivudine- nevirapine)</b>	7	BF240717	R7, 8	0,21231	36,1753	85,71
		DA099283	R7, 4 R7, 8	0,838917	34,1929	
		AW503859	R7, 4 R7, 8	0,838917	34,1929	80,95
		D17012	R7, 4 R7, 8	0,838917	34,1929	
<b>ZDV-3TC-EFV (zidovudine- lamivudine- efavirenz)<sup>a</sup></b>	1		R1,72-3	3,44E-03	32,21	90,48
		DA020938	R1,72-2	0,84	34,19	90,48
			R1,73-13	0,01	40,14	85,71
		T34054	R1 (72-5, 73-9)	0,84	34,19	71,43
		DB095531	R1 (72-6, 73-6)	0,21	36,18	76,19
		BF380718		0,84	34,19	71,43
		BX325384				
		BX340654				
		AL581553	R1 (72-1,72-2, 73- 7, 73-11)			
		AL560565		0,21	36,18	76,19
		AL558996				
		BX325200				
BU633649						
BQ019860						

---

AV757747	72-12				71,43
BF184320	73-7	0,84	34,19		90,48

---

1 <sup>a</sup>Only important hits with significant scores are shown