

# Human immunodeficiency virus type 1 reverse transcriptase and protease mutation search engine for queries

ROBERT W. SHAFER<sup>1</sup>, DUANE R. JUNG<sup>1</sup> & BRADLEY J. BETTS<sup>2</sup>

<sup>1</sup>Division of Infectious Diseases, School of Medicine & <sup>2</sup>School of Engineering,  
Stanford University, Stanford, California 94305, USA

Correspondence should be addressed to R.W.S.; email: [rshafer@cmgm.stanford.edu](mailto:rshafer@cmgm.stanford.edu)

Genetic research is increasingly turning to studies of sequence variation in genes encoding proteins of known structure and function. The principal question in these studies is whether sequence variation affects protein structure or function, and, for certain genes, whether sequence variation affects human health. The proliferation of published sequence data and the growth in the number of publications is a boon to this research, but also makes it difficult to keep track of what is known about a gene. The primary sequence databases of the International Nucleic Acid Sequence Data Library (for example, GenBank) provide powerful sequence-similarity search tools that help researchers deduce the functions of newly identified proteins<sup>1-3</sup>. However, they do not contain the annotation required to map sequence variation within a single gene, or to correlate such variation with data about the gene's product or the phenotype

arising from variation in the gene<sup>4-5</sup>.

Sequence variation is particularly relevant to infectious pathogens that mutate in response to antimicrobial therapy. Sequence variations in human immunodeficiency virus (HIV)-1 reverse transcriptase (RT) and protease, the molecular targets of anti-retroviral drug therapy, are prime examples of genes in which sequence variation has both biological and medical implications. Although HIV-1-infected individuals with drug-susceptible HIV-1 isolates experience substantial reductions in morbidity and mortality with appropriate anti-retroviral drug therapy, individuals infected with drug-resistant isolates generally do not respond to drug therapy<sup>6</sup>.

There are more than 10,000 published HIV-1 RT and protease sequences in GenBank and more than 2,000 additional amino-acid sequences published in biology and medical journals.

Genetic analysis of HIV-1 isolates has demonstrated at least 9 different group M (main) subtypes, differing from one another by 10–30%, and several highly divergent group O and group N outlier sequences<sup>7</sup>. Sequences of HIV-1 isolates collected from around the world provide evidence of naturally occurring genetic polymorphisms; sequences of isolates from patients receiving anti-HIV drug therapy indicate genetic changes selected under anti-retroviral drug pressure that may be associated with drug resistance. The variety and variability of these changes mean that no single researcher can keep track of the consequences of even the more common changes. The combination of this complexity and the lack of annotation in public databases means that most of these published data on repetitively sequenced genes are essentially lost to the research community.

Secondary databases with curated sequences are required to explicate the experimental results stored in the large public primary sequence databases. For example, a publicly available database at the Los Alamos

National Laboratory provides HIV sequence annotation and analysis, and catalogs HIV sequence variation throughout the viral genome, particularly in untreated HIV-1-infected individuals (<http://hiv-web.lanl.gov>)<sup>7</sup>. To allow researchers to analyze new HIV-1 RT and protease sequences arising in patients receiving anti-retroviral drug therapy, we have developed an online search

Table 1 Amino-acid motifs used to align HIV-1 protease and reverse transcriptase sequences

Motif <sup>a</sup>	Positions	Function <sup>b</sup>	% Conserved in HIV-1 <sup>c</sup>	% Conserved in HIV-2
Protease				
PQ(I V)T	1–4	Dimer interface	>99	0
DTG	25–27	Active site	100	100
GCTLN	94–98	Dimer interface	>99	0
RT				
PISP	1–4	Part of protease/ RT cleavage site	>99	0
QWPL(T S)	23–27	β1–αA turn	>99	100
D(V I)GDA	110–114	Active site	>99	100
QY(·)DD(L I)	182–187	Active site	>99	100
WMG(Y F)	229–232	β12–β13 turn	>99	100

<sup>a</sup>I|V, "I" means that either I or V is acceptable at that position; (·), any amino acid is acceptable at the position. <sup>b</sup>DTG constitutes the active site of the protease. The three catalytic aspartates (D) at codons 110, 185 and 186 constitute the active site of the RT. <sup>c</sup>Each motif not only is highly conserved among HIV-1 RT and protease, but also is not found at any other position within these two genes.

Table 2 Mutation frequency data sets used by HIV-SEQ

Gene	Subtype	Treatment	Number of isolates <sup>a</sup>	Number of references <sup>a</sup>
Protease	B	None <sup>a</sup>	436	55
Protease	Non-B	None <sup>a</sup>	163	39
Protease	B	≥1 protease inhibitor	471	32
Protease	B	≥3 protease inhibitors	80	11
RT	B	None <sup>b</sup>	136	42
RT	Non-B	None <sup>b</sup>	117	31
RT	B	≥1 nucleoside RT inhibitors <sup>c</sup>	381	42
RT	B	≥4 nucleoside RT inhibitors <sup>c</sup>	90	14
RT	B	Non-nucleoside RT inhibitors <sup>d</sup>	121	11

<sup>a</sup>These numbers are based on data through 3/1/2000. The most recent figures can be found in the 'Release Notes' section (<http://hivdb.stanford.edu/hiv/programs.htm#ReleaseNotes>). <sup>b</sup>Individuals may have received an RT inhibitor. <sup>c</sup>Individuals in these categories did not receive nonnucleoside RT inhibitors. <sup>d</sup>Individuals in this category usually also received nucleoside RT inhibitors.

## Methods

**User interface.** Users submit single sequences or sets of multiple sequences encompassing HIV-1 protease and/or RT by typing nucleotides into a text box or by uploading a nucleotide text file. The program output includes the amino-acid translation; a list of amino-acid differences (called mutations here) between the submitted sequence and the HIV-1 subtype B consensus sequence<sup>7</sup>; a list of problematic positions indicative of poor sequence quality; a comparison of the submitted sequence to a set of reference sequences of known subtype<sup>9</sup>; data on the frequency with which each mutation occurs in individuals according to HIV-1 subtype and type of drug therapy; background data on potential drug-resistance mutations; and hyperlinks to the Medline references and GenBank entries associated with published sequences. Submitted sequences are not added to the database or stored on the server.

**Sequence translation and alignment.** The program identifies the submitted sequence's correct reading frame by determining which reading frame contains conserved amino-acid motifs (Table 1). The presence of motifs in more than one reading frame indicates a reading frame shift. When this occurs, the program locates, records and removes such shifts by applying an optimal sequence alignment algorithm<sup>10</sup> to the region between the closest motifs in different reading frames. Although reading frame shifts in HIV-1 RT and protease have not been reported, they may conceivably result from defective virions or sequencing errors.

The program then uses the position of conserved amino-acid motifs within the submitted sequence to infer the starting positions of the protease and RT. Amino-acid insertions or deletions will cause the distances between conserved motifs to be different from the inter-motif distances in the consensus reference sequence. When an amino-acid insertion or deletion is detected, the program locates, records and removes it using an optimal sequence alignment algorithm. RT insertions and deletions have been reported in 1–2% of isolates from individuals who have received prolonged treatment with multiple nucleoside RT inhibitors. The occurrence of reading frame shifts and amino-acid insertions or deletions is reported at the beginning of the program's output.

**Representation of amino acid mixtures.** Because of its high mutation rate, HIV-1 exists within an individual as a mixture of genetically distinguishable variants. Sequencing is often done using uncloned, PCR-amplified cDNA (population-based sequencing), and therefore nucleotide mixtures are often detected in clinical plasma virus samples. Such mixtures are represented using the nucleotide ambiguity code (for example, M = A/C, R = A/G and W = A/T). HIV-SEQ translates nucleotide triplets containing ambiguities into each of the possible amino acids they encode. For example, the nucleotide triplet WMC is translated to NTYS (N for AAC, T for ACC, Y for TAC and S for TCC).

program coupled to an extensive database of HIV-1 RT and protease mutations. The program, called HIV-SEQ (for 'HIV mutation search engine for queries') compares an HIV-1 sequence submitted by the researcher to a consensus reference sequence. Differences between the sequences are used as query parameters for 'interrogating' the database<sup>8</sup>, which links HIV-1 RT and protease sequence variation to the anti-retroviral drug treatment histories of patients from whom sequenced virus isolates were obtained. The program is available at <http://hivdb.stanford.edu/hiv/programs.htm>.

**Problematic positions.** The program output includes a list of problematic positions in each sequence. These positions include those coding for stop codons, those containing mutations at highly conserved sites and those containing highly ambiguous nucleotides: N (A/C/G/T), B (C/G/T), D (A/G/T), H (A/C/T) and V (A/C/G). Sequence quality is likely to be inversely correlated with the number of these problematic positions.

**Mutation frequency data.** Each sample mutation is used as a query parameter to 'interrogate' the HIV RT and Protease Sequence Database (Fig. 1). Within the database, mutation frequency tables contain data on the frequency with which each mutation occurs in different categories of HIV-1 isolates. For protease isolates, these categories include subtype B protease isolates from 'untreated' individuals (those who have not received a protease inhibitor, but may have received an RT inhibitor); non-subtype B protease isolates from untreated individuals; isolates from individuals who have received at least one protease inhibitor; and isolates from individuals who have received at least three protease inhibitors.

For RT isolates, the categories include subtype B RT isolates from 'untreated' individuals (those who have not received an RT inhibitor); non-subtype B isolates from untreated individuals; isolates from individuals who have received at least one nucleoside RT inhibitor, but have not received a non-nucleoside RT inhibitor; isolates from individuals who have received at least four nucleoside RT inhibitors, but have not received a non-nucleoside RT inhibitor; and isolates from individuals who have received a non-nucleoside RT inhibitor. The numbers of individuals and literature references associated with each category as of 1 March, 2000 are shown in Table 2.

To minimize reporting bias, the mutation frequency tables contain one sequence per individual. For individuals with multiple isolates, the tables include the earliest isolate from untreated persons and the latest isolate from persons receiving anti-retroviral therapy. To exclude technical sequencing errors and cases of circulating virus containing unusual variants, the tables include only mutations reported in at least two isolates and mutations present as the main form whenever multiple clones from the same isolate were sequenced.

**Background mutation data.** Identified mutations are further annotated with background data that include brief summaries of associations between mutations and drug resistance based on *in vitro* susceptibility test results obtained from recent literature reviews<sup>13,14</sup> and scientific meetings. Background data are available for mutations at 27 of the 99 protease positions and 34 of the 560 RT positions.

**Application deployment.** HIV-SEQ is coded using Perl and CGI programming and runs on a Windows NT operating system<sup>15</sup>. The Win32::ODBC Perl module links the program to tables in the HIV RT and Protease Sequence Database.

## Discussion

HIV-SEQ demonstrates that the body of published sequence data on a gene can be made available in real time to researchers sequencing new isolates of that gene. The ability to examine new sequences in the context of published sequence data has two main advantages. Unusual sequence results can be detected, allowing the person sequencing the gene to recheck the primary sequence output (for example, electropherogram) while it is on his/her desktop<sup>16</sup>. Moreover, unexpected associations between sequences or isolates can be discovered, because the person ana-

