

## Numbering Positions in HIV

# Numbering Positions in HIV Relative to HXB2CG

Bette T. Korber,<sup>1</sup> Brian T. Foley,<sup>1</sup> Carla L. Kuiken,<sup>1</sup> Satish K. Pillai,<sup>1</sup> and Joseph G. Sodroski<sup>2</sup>

<sup>1</sup> *Theoretical Biology and Biophysics, Group T-10, MS K710, Los Alamos National Laboratory, Los Alamos, NM 87545;*

<sup>2</sup> *Howard Hughes Medical Institute, Columbia University, New York, NY 10032*

In this section we present a simple numbering scheme to facilitate the identification of the position number or precise location of interest in HIV DNA or proteins.

Inconsistent and inaccurate numbering of locations in HIV DNA and protein sequences is a serious problem in the HIV literature. Therefore we decided to provide a practical guide to help circumvent these problems in the future, and to attempt to bring a common language into discussions in the field. We present a clearly numbered set of proteins, and the full length genome, for HIV HXB2, GenBank accession number K03455. HIV HXB2 is also known as: HXBc2, for HXB clone 2; HXB2R, in the Los Alamos HIV database, with the R for revised, as it was slightly revised relative to the original HXB2 sequence; and HXB2CG in GenBank, for HXB2 complete genome. Our web site has an interactive program to further facilitate obtaining position numbers relative to HXB2CG (<http://hiv-web.lanl.gov/NUM-HXB2/HXB2.MAIN.html>).

HXB2 was selected as the prototype, because this virus is the most commonly used reference strain for many different kinds of functional studies. Importantly, all of the envelope structural data published to date translates residue numbers into the HXB2 numbering scheme. Now that a core HIV-1 gp120 structure is solved (for review, see Wyatt et al, this compendium), it has become apparent that conservation in core sequences, especially in hydrophobic interior domains, exists to preserve similar folding in gp120 variants. As the envelope protein is riddled with insertions and deletions, it is particularly problematic for numbering. The current system, of sequentially numbering proteins from any strain, lacks a common way to refer to specific locations in a protein. We propose the following system to circumvent this problem:

- 1) ***Case of insertion in sequence relative to HXB2CG.*** Use residue number/alphabet (e.g., 131a, 131b, 131c, etc.) to refer to residues in variable regions that are “extra” compared to what HXB2 has. A similar scheme has

been used for immunoglobulin complementarity-determining region (CDR) loops (see Lucas et al., *J Immunol* 1998 **161**:3776–80 (1998) for an example).

Example: If the region under study is LLITRDGGSNRSEPEVEIFRP of ENVB, gp120,

452 | 465 | 470 HXB2 amino acid position from start of gp160  
LLITRDGGSNNSNES-EIFRP HXB2  
LLITRDGGSNRSEPEVEIFRP ENVB

one could refer to it as corresponding to HXB2 gp160 position numbers 452–470 with a two base insertion (465a = E and 465b = V)

- 2) ***Case of deletion in sequence relative to HXB2.*** Indicate the deleted residues.

Example: If the region under study is LLITRDGGNN of 92RW020.5,

452 | 463 HXB2 amino acid position from start of gp120  
LLITRDGGNNNN HXB2  
LLITRDGG..NN 92RW020.5

one refers to it as corresponding to HXB2 gp160 position numbers 452–463 with a two base deletion at positions 460–461. We suggest using the annotation 452-463(del 460-461) to make this explicit.

The sequential numbering relative to either 92RW020.5 or ENVB could also be provided in the above two examples, but the HXB2 numbering should also be provided as a reference.

The benefit of this numbering strategy is that, for example, aspartate 368, which is involved in CD4 binding, or gp160 368 D, means the same thing to everyone working on envelope glycoproteins, regardless of the reference strain they used in their particular studies.

Also, when working with a short functional domain, epitope, or primer, researchers should publish the precise amino acid or nucleotide string that they are working with, as well as the HXB2 numbered positions, to ensure that there is no confusion (for example, write out ENVB LLITRDGGSNRSEPEVEIFRP as well as give the boundary position numbers).

## Numbering Positions in HIV

### HXB2 Amino Acid Sequence Numbering:

---

#### Gag Pr55 Gag precursor (Assembly)

MGARASVVLSG GELDRWEKIR LRPGGKKKYK LKHTIWASRE LERFAVNPGI LETSEGRQI LGQLQPSLT GSEELRSILN TVATLYCVHQ RIEKDTKEA 100  
LDKTEEEQNK SKKKAQQAA DTGHSNQVSQ NYPIVQNIQG QMVHOAISPR TLNAWVKVVE EKAFSPEVIP MFSALSEGAT PQDINTMLNT VGGHOAMQM 200  
LKETINEEAA EMDRVHPVHA GPIAPGQMR E PRGSDIAGTT STIQEQIGMM TNPPPIPVG E IYKRWIIILGL NKIVRMSPT SILDIRQGPK EPFRDYVDRF 300  
YKTLRRAEQAS QEVKNWMTET LIVQNANPDC KTIKALGP A TLEEMMTAC QGVGGPGHKA RVLAEAMSQV TNSATIMMQR GNFRNQRKTV KCFNCGKEGH 400  
TARNCRAPRK KGCWKCGKEG HQMKDCTERQ ANFLGKIWPS YKGRPGFLQ SRPEPTAPPE ESFRSGVETT TPPQKQEPID KELYPLTSLR SILFGNDPSSQ 500

#### Gag p17 Matrix

MGARASVVLSG GELDRWEKIR LRPGGKKKYK LKHTIWASRE LERFAVNPGI LETSEGRQI LGQLQPSLT GSEELRSILN TVATLYCVHQ RIEKDTKEA 100  
LDKTEEEQNK SKKKAQQAA DTGHSNQVSQ NY 132

#### Gag p24 Capsid

PIVQNTIQQM VHQAISPR TL NAWKVVEEK AFSPEVIPMF SALSEGATPQ DLNTMLNTVG GHQAAMQMLK ETINEEAAEW DRVHPVHAGP IAPGQMRP 100  
GSDIAGTTST LQEQQIGWMTN NPPIPVGELY KRWWILGLINK IVRMYSPTSI LDTRQGPKEP FRDVYDRFYK TLRAEQASQE VKNWMTETILL VQNAINDCCT 200  
ILKALGPAAT LEEMMTACQG VGGPGHKARV L 231

#### Gag p2

AEAMSOVTNS ATIM 14

#### Gag p7 Nucleocapsid

MQRGNFNRNQR KIVKCFNCGK EGH TARNCR A PRKKGCWKG C KEGH QMKDCT ERQAN 55

#### Gag p1

FLGKIWP SYK GRPGNF

#### Gag p6

LQSRPEPTAP PEEFSRSGVE TTTPPKQEP IDKELYPLTS LRSIFGN DPS SQ 52

## Numbering Positions in HIV

**Pol polyprotein:**

FREDDLAFLQ	GKAREFSSEQ	TRANSPTRE	LQWGRDNNS	PSEAGADRGQ	TVSFNFQVTT	LWQRPLVTIK	IGGQKEALL	DTGADDTIVL	EMSLPGRWKP	100
KMIGGIGGFI	KVRQYDQILI	EICGHKAIGT	VLVGPTPVNI	IGRNLITQIG	CTLNFPISPI	ETVPVKLKG	MDGPVKVQWP	LTEEKIKALV	EICTEMEKEG	200
KISKIGPENP	INTPVFAIKK	KDSTKWRKLV	DFREIINKRTQ	DFWEVQLGIP	HAGLIKKKKS	VTVLVDVGDAY	FSVPLDEDFR	KYTAFITPSI	NNETPGIRYQ	300
YNVLQPQGWKG	SPAIFQSSMT	KILEPFRKQN	PDIVIYQYM	DLYVGSDLEI	GQHRTKIEEL	RQHLLRWGLT	TPDKHHQKEP	PFLWMGYELH	PDKWTWQPTV	400
LPEKDSWTVN	DIQKLVKGKLN	WASQIYPGIK	VRLQCKLIRG	TKALTEVIPL	TEAEELELAE	NREILKEPVH	GVIYDPSKDL	IAEIQOKQGQ	QWTYQTYQEP	500
FKNLKTGKYA	RMRGAHTNDV	KOLTEAVQKI	TTESIVIWGK	TPKFKLIPIQK	ETWETMWTEY	WQATWIPEWE	FVNTPPLVKL	WYQLEKEPIV	GAETFTYVDGA	600
ANRETKLGKA	GYVTNRGRQK	VVTIITDTNQ	KTELQAIYIA	LQDSGLEVENI	VTDSQYALGI	IQAOPDQSES	ELVNOITEQL	IKKEKVYLAW	VPAHKGGNN	700
EQVDKLVSAG	IRKVLFLDGI	DKAODEHEKY	HSNWARAMASD	FNLPDVVAKE	IVASCDCQCL	KGEAMHGQD	CSPGIWQLDCL	THLEGKVLL	AHVHASGYIE	800
AEVIPAETGQ	ETAYFLLKLA	GRWPVKTIHT	DNGSNFTGAT	VRAACWVAGI	KQEFGIPYNP	QSQGVVESMN	KELKKIIQGV	RDQAELHHTA	VQMAVFHNF	900
KRKGGIGGYS	AGERIVDIA	TDIQTKELQK	QITKIQNFRV	YYRDSRNPLW	KGPAKLLWKG	EGAVVIQDNS	DIKVVPRRKA	KIIRDYQKQM	AGDDCVASRQ	1000
DEP						1003				

**Pol p10 Protease**

PQVTLWQRPL VTIKIGGQLK EALLDTGADD TVLEEMSLEP RWPKPMIGGI GGFIKVROQD QILIEIGHK AIGTVLVGPT PVNIIGRNLL TQIGCTLNF 99

**Pol p66 Reverse Transcriptase (RT/RNAse)**

PISPIETV PV	KLKPGMDGPK	VKQWPLTEEK	IKAVALCITE	MEKEGKISKI	GPNPYNTPV	FAIKKKDSTK	WRKLWDFREL	NKRTQDFMEV	QLGIPHPAGL	100
KKKKSVTV LD	VGDAYF SVPL	DEDFFRKYTA F	TIPSINNETP	GIRYQVNVL P	QGNKGSPAIF	QSSMTKILEP	FRKQNPDIVI	YQYMDDLVYG	SDLEI1QHRT	200
KIEELRQHLL	RWGLTT PDKK	HQKEPPFLWM	GYELHPDKWT	VOPTIVIPEKD	SWTVNDIQKL	VGKLNWASQI	YPGIKVRLC	KLURGTKALT	EVIPLTEEAE	300
LELAENRE IL	KEPVHGVYYD	PSKDLIAEIQ	KQGQGQWTYQ	IYQEPFKNLK	TGKYARMRGA	HTNDVKQLT	AVQKITTESI	VIWGKTPFK	LPIQKETWET	400
WWTEYWQATW	IPWEWFVNTP	PLVWKWYQLE	KEPIVGAETF	YVDGAANRET	KLGKAGYVTN	RGRQKVVTLT	DTTNQKTELQ	AIYLALQDSG	LEVNIIVDSQ	500
YALGIIQACP	DQSESELV NQ	IEQLIKKEK	VYLAWVPAHK	GIGGNEQV DK	LVSA GIKV L					560

**Pol p51 RT**

PISPIETV PV	KLKPGMDGPK	VKQWPLTEEK	IKAVALCITE	MEKEGKISKI	GPNPYNTPV	FAIKKKDSTK	WRKLWDFREL	NKRTQDFMEV	QLGIPHPAGL	100
KKKKSVTV LD	VGDAYF SVPL	DEDFFRKYTA F	TIPSINNETP	GIRYQVNVL P	QGNKGSPAIF	QSSMTKILEP	FRKQNPDIVI	YQYMDDLVYG	SDLEI1QHRT	200
KIEELRQHLL	RWGLTT PDKK	HQKEPPFLWM	GYELHPDKWT	VOPTIVIPEKD	SWTVNDIQKL	VGKLNWASQI	YPGIKVRLC	KLURGTKALT	EVIPLTEEAE	300
LELAENRE IL	KEPVHGVYYD	PSKDLIAEIQ	KQGQGQWTYQ	IYQEPFKNLK	TGKYARMRGA	HTNDVKQLT	AVQKITTESI	VIWGKTPFK	LPIQKETWET	400
WWTEYWQATW	IPWEWFVNTP	PLVWKWYQLE	KEPIVGAETF	YVDGAANRET	KLGKAGYVTN	RGRQKVVTLT	DTTNQKTELQ	AIYLALQDSG	LEVNIIVDSQ	500
YALGIIQACP	DQSESELV NQ	IEQLIKKEK	VYLAWVPAHK	GIGGNEQV DK	LVSA GIKV L					560

**Pol p15 RNase**

YV DGAANRET KLGKAGYVTN RGRQKVVTLT DTINQKTELQ A IYLA LQDSG LEVNIVTDSQ YALGIIQACP DQSESELV NQ I IEQLIKKEK VYLAWVPAHK 100  
GIGGNEQV DK LVSA GIKV L 120

**Pol p31 Integrase**

FLDGIDKA QD	EHEKYHSNWR	AMASDFNLPP	VWAKEIVASC	DKCQIKGEAM	HGQVDCSPGI	WQDCTHLEG	KVILVAVHVA	SGYIEA EVIP	AETGQETAYF	100
LLKLAGRWPV	KTIHTDNGSN	FTGATVRAAC	W MAGIKO EFG	IPXNPOSQGV	VE SMNKELKK	II GQVRDQAE	HLK TAVQMAV	FTHNFKRKGG	IGGYSAGERI	200
VDIIATDQT	KELQKQITKI	QNF R VVYRDS	RNPLWKGPAK	LIWKGEGAVV	I QDN SDIKVV	PRRKAKIIRD	Y GKQWAGDDC	VAS QDED		288

## Numbering Positions in HIV

**Vif**

MENRWQVMIV WQVDRMRIRT WKSLVKHHMY VSGKARGWY RHYESPHPR ISSEVHILG DARLVITYW GLHTGERDWL LGQGVSIERW KKRYSIQVDP 100  
ELADQLIHL YFDCFSDSAI RKALLGHIVS PRCEYQAGHN KVGSLOYLAL AALITPKKIK PPLPSVTKL EDRWNKPKQT KGHRGSHTMN GH 192

**Vpr**

MEQAPEDQGP QREPHNEWTL ELLEELKNEA VRHFPRIMIH GIGQHIVETV GTIWAGVEAI IRILQQLIFI HFRIGCRHSR IGVIRQRAR NGASRS 96  
MEPVDPRLEP WKHPGSQPKT ACTNCYCKC CFHCQVCFIT KALGISYGRK KRRRRRAHQ NSQTHQASLS KOPTSQPRGD PTGPKE\$KKK VERETETDPF 100  
D

**Tat**

( premature HXB2 stop codon indicated by \$ )  
MAJ3RSGDSDE ELIRTVRLIK LLYQSNPPPN PEGTRQARRN RRRWREROR QIHSISERIL GTYLGRSAEP VPLQLPPLER LTLDQNEDCG TSGIQQVGSP 100  
QILVESPTVL ESGTE 116

**Vpu**

( defective start codon )  
TOPPIPTAVI ALVVAIIIAI VVWSIVIEY RTKLRQRKID RLIDLIERA EDSGNESEGE ISALVEMGVE MGHHAPWDVD DL 82

**Envelope (Env) gp160**

Env signal peptide |

MRVKEKYQHL WRWGWRWGM LIGMIMICSA TEKLWVTVY GVPVWKEATT TLFCASDAKA YDTEVHNWVA THACVPTDPN PQEVVLVNT ENENMKNDM 100  
VEQMHEDIIS LWDQSLKPCV KLTPICVSLK CTDLKNDTNT NSSSGRMIVE KGEIKNCFSN ISTSIRGKVQ KEYAFFYKLD IIPIDNDT'S YKLTSCNTSV 200  
ITZACPVKSF EPIPIHYCAP AGFAALKCNN KTFNGTGPCCT NVSTVQCTHG IRPVVSTQLL INGSLAEEEV VIRSVNFTDN AKTIIIVQINT SVEINCTRPN 300  
NNTRKIRIQL RGPGRAFTI GKIGNMRQAH CNISRAKWNIN TLKQIAASKLR EQFGNNKTII FKQSSGGDPE IVTHSFNCGG EFFYCNSTQI FNSTWNNSTW 400  
STEGSNNTEG SDTITLPCRI KQINNMWQKV GKAMYAPPIS GOIRCSSNIT GLLTRDGGN SNNESEIFRP GGGDMRDNWR SELYKYYKVK IEPLGVAPTK 500  
> 9P41 start  
AKRIVVQREK RAVGIGALFL GFLGAAGASTM GAASMLITVQ ARQIISGTVQ QDNLLRATIE AQQHLLQLTW WGIKQLQART LAVERYLKDQ QLIGINGCSG 600  
KLICTTAVPW NASWSNKSLE QIWNHTTWME WDREINNYTS LIHSLTEESQ NOQEKNQEL LEIWKWASLW NWFNITNWLW YIKLFIIMVG GIVGLRIVFA 700  
VLSIVNRVRQ GYSPLSFTQH LPTPRGPDRP EGIEEgger DRDRSTRILVN GSALIWWDDL RSLCLFSYHR LRDILLIVTR IVEILLRRGN EALKYWNLL 800  
QWMSQELKNS AVSLLNATAI AVAEGTDRV1 FVWQGACRAI RHPRRIRQG LERILL 856

**Nef** ( premature HXB2 stop codon indicated by \$ )

MGGKWSKSSV IGMPPTVTERM RRAEPAADRV GAASRDLEKH GAITSSNTAA TNAACAWLEA QEEEFVGFPV TPQVPLRPMT YKAAVDLSHF LKEKGLEG 100  
IHSQRQDIL DLWIYHTQY FPD\$QNYTPG PGVRYPLTFG WCYKLVLPVEP DKTEEANKGE NTSLLHPVSL HGMDDPEREV LEWRFDSRLIA FHIVARELHP 200  
EYFKNC 206

## Numbering Positions in HIV

### HXB2 Nucleotide Sequence Numbering:

> 5' LTR U3 region start  
tggaaaggct aattcactcc caacgaagac aagatatcc ttatctgtgg atctaccaca cacaaggcta cttccctgtat tagcagaact acacaccagg  
gccaggatc agatatcac tgacctttgg atgggtgtac aagcttagtac cagttgagcc agagaagtttta gaagaaggcca acaaaggaga gaacaccagg 100  
ttgttacacc ctgtgagct gcatggaaatg gatgaccgg agagagaagt gtttagatggg aggtttgaca gcccgcctagc atttcatacc atggcccgag 200  
auctgcattc ggagttacttc aagaactgtc gacatcgagc ttgttacaag ggactttccg ctgggactt tccaggagg cgtggctgg gccccactgg 300  
ggagtggcga gccctcagat cctgcataata agcagactgtctttgttgcctgt actgggtata tctggtttaga ccagatctga gcctgggagc tatctggcta 400

5' LTR U3 region end \ 5' LTR R repeat start  
ggagtggcga gccctcagat cctgcataata agcagactgtctttgttgcctgt actgggtata tctggtttaga ccagatctga gcctgggagc tatctggcta 500

5' LTR R 5' LTR U5  
repeat end \ region start  
actagggAAC ccactgttta agcctcaata aagcttgct tgagtgttg tgccctgt ttgtgtgact ctggtaacta gagatccctc 600

5' LTR U5 region end <  
agaccctttt agtcagtgtg gaaaatctct agcagtggc cccgaaacagg gacctgaaag cgaaggaa accagaggag ctctctcgac gcaggactcg 700

gcttgcgtaa ggcgcacgg caagggcga gggggcggcga ctgggtgatg cggccaaat tttgactagc ggaggctaga aggagagaga tgggtgcgg 800  
aegctcagta ttaagcgggg gagaatttga tcgatggaa aaaattcgt taaggccagg gggaaagaaa aaatataaat taaaacatat aatgtggca 900  
agcaggaggc tagaacgtt ctcggctgt tagaaacatc agaaggctgt agacaatac tggacagct acaaccatcc ctccagacag 1000  
gatcagaaga acttagatca ttatataata cagtagcaac cctcttattgt gtgcataaa ggatagagat aaaagacacc aaggaagct tagacaagat 1100

> Gag p17 start  
agaggaagag caaaacaaa gtaaaaaaa agcacacgaa gcaagcgtg accacaggaca cagaatctag stcggccaa attaccat atgtgcagaac 1200  
atccaggggc aaatgtaca tcagccata tacatggaa cttaaatgc atgggtaaa gtagttagaa agaaggctt cagcccgaaa gtatcacca 1300  
tgttttcagc attatcgaa ggagccaccc cacaagatt aaacaccatg ctaaacacag tggggaca tcaagcagcc atgcaaatgt taaaagagac 1400  
catcaatcgag gaagctcgag aatggatag aytgcatccca gtgcatgcag ggcctattgc accaggccg atgagagaac caagggaaatg tgacatcgaa 1500  
ggaactacta gtaccctca ggaacaataa gatggatgaa caataatcc acctatccca gttaggaaaa ttataaaag atggataatc ctgggataa 1600  
ataaaatagt aagaatgtat agcctacca gcatcttggaa cataagacaa ggaccaaaagg aaccctttag agactatgtt gacccgttct ataaaactct 1700  
aaagccgag caagttcac agggtaaa aaatggatg acagaaacctt tgggttcaaa aaatgcgaa ccagattgtt agactatgtt aaaaactctg 1800

Gag p24 Capsid end \ Gag p2 start  
ggaccagcgg ctacactaga agaaatgtg acagcatgtc agggagtagg aggaccggc cataaggaa gagtttggc tgaagcaatgg agccaagttaa 1900

## Numbering Positions in HIV

Gag p2 end \ / Gag p7 Nucleocapsid start  
 caaattcagc taccataatg atgcagagag gcaattttag gaaccaaaga aagatgtta agtgttcaa ttggccaaa gaagggcaca cagccgaaa 2000

ribosome -1 slip Gag to Gag-Pol

---

Gag p7 nucleocapsid end \ /Gag p1 start  
 Pol start >  
 Pol protease end \ / Pol p66 and p51 RT start  
 Gag p6 end <

ctccccctca	gaagcagag	ccgatagaca	aggaactgtta	tcctttact	tcctcgagt	cactcttgg	caacgacccc	tgcacacaat	aaagataggg	2300
gggcaactaa	aggaagctt	attagataca	ggagcagatg	atacagtatt	agaagaatg	agtttgccag	gaagtgaa	acccaaatgg	ataggggaa	2400
tggaggtt	tatcaaagta	agacagatg	atcagatact	catagaatc	tgggacata	aagtatagg	tacagtata	gttggaccta	cacctgtcaa	2500

> Pol protease start

Pol p51 end p66 RT continue \ / Pol p15 RNase H start  
 aatccccctc

ccttagtggaa	attatggta	cagttagaga	aagaacccat	agtaggagca	gaaaccttct	atgtagatgg	ggcagctaac	agggagacta	3900	
atttaggaaa	agcaggatat	gttactata	gaggaagaca	aaaagtgtc	accctaactg	acacaacaa	tcaagagact	gagttacaag	caattttatct	4000
actttgcag	gattcggat	tagaagtaaa	catagtaaca	gactcacaat	atgcatttag	aatcattca	gcacacaa	atcaaagtgt	atcagagttt	4100
gtcaatcaa	taatagagca	aggaaaaagg	tctatctggc	atgggttaca	gcacacaa	gaatggagg	aatatgaacaa	gtagatataat		4200

## Numbering Positions in HIV

Pol RNASE H, p66 RT end \ Pol p31 Integrase start  
 tagtcagtgc tggaaatcagg aaagtactat ttttagatgg aatagataag gccaaagatg aacatgagaa atatcacagt aattggagag caatggctag 4300  
 tgatttaac ctggccactg tagtagcaa agaaatatgt a gccagctgtg ataaatgtca gctaaagga gaagccatgc atggacaagt agactgttgt 4400  
 ccaggaatat ggcactataga ttgtacacat ttagaaggaa aagttaatctt ggttgcatgtt catgtggccca gtggatata agaaggcagaa gtttcccg 4500  
 cagaacagg gcaggaaca gcatattttc ttttaaaattt agcggaaaga tggccagtaa aaacaataca tactgacaat ggcagcaatt tcaccgg 4600  
 tacggtagg gccgcctgtt ggtgggggg aatcaaggcag gaatttggaa ttccctacaa tccccaaatgtt caaggatgt tagaatctat gaataaaagaa 4700  
 ttaaagaaa ttatggaca ggttaagatg caggctgaa acatctaagc agcgtacaa atggcagtat tcatccacaa ttttaaaaga aaagggg 4800  
 tggggggta cagtgcagg gaaagaatag tagacataat agcaacagc atacaacta aagaattaca aaaacaaattt acaaaaatttccg 4900  
 ggtttattac agggacagca gaaatccact ttggaaagga ccagcaaagc tcctctggaa aggtgaagg gcagtagtaa tacaagataa tagtgacata 5000

> Vif start Pol, p31 Integrase end <  
 aaagttagtgc caagaagaaa agcaaaagtc attaggatt atggaaaaca gatggcagggt gatgatgtg tggcaagtag acaggatgg gatttagaaca 5100  
 tggaaaagttagtacaa ccatatgtat gtttcaggaa aagcttaggg atgggtttat agacatctact atgaaagccc tcatccaaga ataagtcc 5200  
 aagtacacat cccactagggtatgtcataat aacatattgg ggtctgcata caggagaag agactggcat tgggtgcagg gagtctccat 5300  
 agaatggagg aaaagagat atagcacaca aatagacccctt gaacttagcag accaactaat tcatctgtat tacttgact gtttttcaga ctctgcata 5400  
 agaaaggcct tattaggaca catagttac ccttaggtgtt aatatacaggc aggacataac aaggttagat ctctacaata cttggcacta gcagcaata 5500

taacacaaa aaagataaag ccactttgc ctagtgttac gaaactgaca gaggatagat ggaacaagcc ccagaagacc aagggccaca gagggagcca 5600

Vif end <  
 cacaatgaat ggcactataga gcttttagag gagcttaaga atgaagctgt tagacatttt ccttaggattt ggctccatgg ctttagggcaa catatctatg 5700  
 aaacttatgg ggtactctgg gcaggagttgg aagccataat aagaattctgt caacaactgtc tttttatcca ttttcagaat tgggtgtcga catagragaa 5800

Tat start > Vpr end <  
 taggcgttac tcgacagagg agagcaagaa atggagccat tagatccatg actagagccc tggaaagcatc caggaagtca gcctaaact gcttgatcca 5900

Rev start >  
 attgcttatgt taaaatgt tgatccatt gccaagtttgg tttccataaca aagccattag gcatctctca tggcaggaaag aagcggagac agcggaaag 6000

Tat, Rev exon end \Tat, Rev intron > Vpu start (defective ACG start codon)  
 agctcatcag aacagtcaaa ctcatcaagc tctcttatca aagcagtaag tagtacatgt aacgcacact ataccaatag tagcaatagt agcatttagta 6100  
 gtagcaataa taatagcaat agttgtgtgg tccatagtaa tcatatagaata tagaaataa ttaagacaa gaaaataga caggtaattt gatagactaa 6200

> Env gp160 start, signal peptide  
 tagaaagagc agaagacagt ggcataatgaga gtgaaggaga aatatcagca ctgtgtggaa tgggggtgaa gatgggcac catgtctactt gggatgttga 6300

## Numbering Positions in HIV

Vpu end, signal peptide end

<  
tgatctgttag tgctacagaa aatagtggg tcacagtcta ttatgggta cctgtgtgga aggaagca caccactcta ttttgtcacat cagatgctaa 6400  
acatatgtat acagaggta ataatgtttg ggcccacat gcctgtgtac ccacagaccc caaccacaa gaagtagtat tggtaaatgt gacagaaaat 6500  
ttAACATGTG gaaaaatga catgttagaa cagatgcattt aggatataat cagtttatgg gatcaagcc taaagccatg tgtaaaatgtt acccaactt 6600  
gtgttagtt aaagtgcact gattgaga aatgataactaa taccatataatc agtagcgaa gaatgataat ggagaaagga gagataaaaaa actgctctt 6700  
caatatcgc acaaggctaa gaggttagt gcagaaagaa tatgcattt ttataaact tgatataata ccaatagata atgataactac cagctataag 6800  
ttgacaagt gtaacactc agtattaca caggctgtc caaggatc cttgagca attccatc atttgtgc cccgctgg tttgcattc 6900  
taaatgtaa taataagac ttcaatggaa caggaccatg tacaatgtc agcacagtc aatgtacaca tggattagg ccagtagtat caactcaact 7000  
gctgttaat ggcagtcg cagaagaaga ggttagtaat agatctgtca atttcacgga caatgtcaa accataatag tacagtgaa cacatctgt 7100  
gaaattaatt gtacaagac caacaacaat acaagaaaaa gaatccgtat ccagagagga ccaggagag cattgtac aataggaaaa ataggaaata 7200  
tgagacaagc acatgttac attagtagag caaatggaa taacacttta aaacagatag ctgcaaaatt aagagacaa tttggaaata ataaaacaat 7300  
aatcttttaag caatcctcag gaggggaccc agaaatgttta acgcacagtt ttaatgtgg agggaaattt ttctactgtta attcaacaca actgtttat 7400  
agtacttggt ttaatgtac ttggagtagt gaaagggtcaa ataacactga aggaagtgtac acaatcccc tccatgcag aataaaacaa attaaaca 7500  
tgtggcagaa agtagaaaaa gtaatgtatc agtggacaa attagatgtt catcaaatat tacaggctg ctatcaaaca gagatgtgg 7600  
taatagcaac aatgagtcg agatcttcg acctgggaga ggagatattga gggacaattg 7700

Env gp120 end \ Env gp41 start

ccattaggag tagaccac caagcaag aagaaagtgg tcagagaga aaaagagca gtggaaatag gagcttggt ccttgggtt ttggagcag 7800  
caggaagcac tatggcga gcctcaatga cgctgacggt acaggcaga caattatgt ctggatagt gcagcagcag aacaatttgcc tgagggtat 7900  
tgaggcgca cagcatctgt tgcaactcac agtctgggg atcaaggcagc tcggagaag aatctgtgt gtggaaatgt acctaataagg tcaacagctc 8000  
ctggggatgt ggggtgtct tggaaactc atttgccaca ctgtgtggc ttggaaatgtt aatgtggata ataaatctt ggaacagat tggaaatcaca 8100  
cgacactggat ggagtggac agagaattt acaatttacaa aagtttataa caactctttaa ttggtatataa aaattattca taatgtatgtt aggaggctt 8200  
attattggaa ttgatataat gggcaattttt ttggaaatttg ttgatataatc caaaatggctt 8300

Tat, Rev intron end \ Tat, Rev exon 2 start  
gtaggtttaa gaatagttt tgctgtactt tctatataatc atagagtttgcagatggat tcaccattat cgtttcagac ccacccccc accccgggg 8400

— Tat premature stop

Tat end <

gaccggacag gcccgaagaa atagaagaag aaggtggaga gagagacaga gacatcca ttcgattatgtt gaaaggatcc ttggcactta tctgggacga 8500  
tctggggaggc ctgtgcctt tcagttacca ccccttgaggaa gacttacttctt tgatgttac gaggatgtt gaaacttctgg gacgggggg gttggaaatcc 8600

Rev end <

ctcaaatatt ggtggaaatctt cctacacgtat tggacttcagg aactaaagaa tagtgcgtt agcttgcgtca atgcacacgc catagcagta gctgagggg 8700

## Numbering Positions in HIV

Env gp41, gp160 end < > Nef start  
cagataggt tataagaaga gtaccaaggag cttgttagagc tattcccac atacctagaa gaataagaca 999cttggaa aggattttgc tataagatgg 8800  
gtggcaagtg gtcaaaaagt agtgtgattg gatggcctac tgtaaggaa agaatgagac gagctgagcc 8900  
agacctggaa aaacatgag caatcacaag tagcaataca gcagcttcca atgctgcttg tcctggcta gaagcacaag 9000  
  
> 3' LTR U3 region  
ccagtcacac ctcaaggtacc tttaagacca atgacttaca aggcagctgt agatcttagc cacttttaa aagaaggg gggactggaa 999ctaattc 9100  
  
actcccbaag aagaacaat atccttgc tttggatctt ccacaccaa ggctacttcc 9200  
tccactgacc ttggatgt gctacaagct agtaccagg tggccagata agatagaaga 9300  
acctgtcatg ggtggatga cccggagaga gaaagtgttag agtggagtt tgacagccgc 9400  
  
Nef end <  
acttcaagaat ctgctgacat cgagcttgc acaaggact ttccgtgg gactttccag ggaggcgtgg 9500  
  
3' LTR U3 region \ / 3' LTR R repeat  
cagatctgc atataaagcag ctgcttttg cctgtactgg gtctctctgg ttagaccga tctgagctg ggagatctt ggtaacttag ggaaccact 9600  
  
3' LTR R repeat \ / 3' LTR U5 region  
gtttaaggct caataaagct tgcttggatg gttcaagta gtgtgtggcc gtctgtgtg tgactctggt aacttagagat ccctcagacc cttttagtca 9700  
  
3' LTR U5 end <  
gtgtggaaaa tctctagca 9719

This numbering was based on previous HIV sequence database annotation, cross-checked with protein structure databases, Tozser et al., *FEBS letters* **281**:77–80 (1991), and R. J. Gorelick and L. E. Henderson, *Human Retroviruses and AIDS 1994*, part III, pages 2–10.