

Numbering Positions in HIV Relative to HXB2CG

Bette T. Korber,¹ Brian T. Foley,¹ Carla L. Kuiken,¹ Satish K. Pillai,¹ and Joseph G. Sodroski²

¹ Theoretical Biology and Biophysics, Group T-10, MS K710, Los Alamos National Laboratory, Los Alamos, NM 87545;

² 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 LLLTRDGGSNRSEPEVEIFRP of ENVB, gp120,

452	465	470	HXB2 amino acid position from start of gp160
LLL	TRDGGSNNES	--EIFRP	HXB2
LLL	TRDGGSNRSEPEVE	IFRP	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 LLLTRDGGNN of 92RW020.5,

452	463	HXB2 amino acid position from start of gp120
LLLTRDGGGSNN		HXB2
LLLTRDGG..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 LLLTRDG-GSNRSEPEVEIFRP as well as give the boundary position numbers).

We intend to change the HIV Immunology Database to this system, through the course of 1999. This year we have made the HXB2 strain the reference strain in our alignments for the sequence compendium, although the WEAU strain remains the reference strain for the immunology compendium in 1998.

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.

Numbering Positions in HIV

HXB2 Amino Acid Sequence Numbering:

Gag Pr55 Gag precursor (Assemblin)

MGARASVLSG GELDRWEKIR LRPGGKKYK LKHIVWASRE LERFAVNGL LETSEGCRQI LGOLOPSLOT GSEELRSLYN TVATLYCVHQ RIEIKDTKEA 100
LDKIEEEQNK SKKKAAQAAA DTGHSNOVSQ NYPIVONIQG OMVHOAISPR TLNAWVKVVE EKAFFSPEVIP MFSALSEGAT PQDLNTMLNT VGGHQAAMQM 200
LKKETINEEAA EWDRVHPVHA GPIAPGOMRE PRGSDLAGTT STLQEQQIGWM TNNPPIPVGEG IYKRWLILGL NKIVRMYSPV SILDIRQGPK EPFRDYVDRF 300
YKTLRAEQAS QEVKNWNTET LLVQNANPDC KTIILKALGPA ATLEEMMTAC QGVGGPGHKHA RVLAFAAMSQV TNSATIMMRQ GNFRNRQKIV KCFNCGKEGH 400
TARNCRAFRK KGGWRCGKGEG HQMKDCTERQ ANFLGIWPS YKGRGPGNFLQ SRPEPTAFAPE ESFRSGVETT TPPQKQEPID KELYPLTSLR SLFGNDPSSQ 500

Gag p17 Matrix

MGARASVLSG GELDRWEKIR LRPGGKKYK LKHIVWASRE LERFAVNGL LETSEGCRQI LGOLOPSLOT GSEELRSLYN TVATLYCVHQ RIEIKDTKEA 100
LDKIEEEQNK SKKKAAQAAA DTGHSNOVSQ NY 132

Gag p24 Capsid

PIVQNIQCM VHQAISPRTL NAWKVYEEK AFSPEVTPMF SALSEGATPQ DLNTMLNTYG GHQAMQMLK ETINNEAAEW DRVHPVHAGP IAPGQMREPR 100
GSDIAGTTST LQEQLGWTN NPPIPVGELY KRWLILGLNK IVRMYSPTSI LDIRQGPREP FRDYVDRFYK TLRAEQASQE VKNWMTETLL VQNANPDCKT 200
ILKALGPAAT LEEMMTACQG VGGPGHKARV L 231

Gag p2
AEAMSQVTNS ATIM 14

Gag p7 Nucleocapsid

MORGNNFRNQR KIVKCFNCGK EGHTARNCRA PRKGKWKCG KEGHQMKDCT ERQAN 55

Gag p1
FLGKIWP SYK GRPGNF 16

Gag p6

LQSRPEPTAP PEESFRSGVETT TTPQKQEP IDKELLYPLTS LRSLLFGNDPS SQ 52

Pol polyprotein:

FREDLAAFLQ GKAREFSSEQ TRANSPTRRE LOVWGRDNNS PSEAGADRGQ TVSFNFPOVT LWQRPLVTIK IGGOLKEALL DTGADDVLE EMLPGRWKP 100
 KMIGGIGGFI KVRQYDILLI EICGHKAIGT VLVGPTPVNI IGRNLLTQIG CTLNFPISP IETVPKVKQWP MDGPVKVKQWP LTEEKIKALV EICTEMEKG 200
 KISKIGPENP YNTPVFAIKK KDSITKWRKLIV DFRELNKRKTQ DFWEVQLGIP HPAGLKKKKSS VTVLDVGDAY FSVPLDDEFR KYAFTIPS1 NNETPGIRYQ 300
 YNVLPQGWKG SPAIEFQSSMT KILLEPRKQN PDIVIYQMD DLYVGSDELQ QHRTKIEEL RQHLLRWGELT PDKKHKQKEP PFLWMGYELH PDKWTQPIV 400
 LPEKDSWTVN DIQKLVYGLN WASQIYPGIK VRQLCKLLRG TKALTEVIPL TEEAEELIAE NREILIKEPVH GVVYDPSKDL IAEIQKGQG QWTYQIYQEP 500
 FKNLKTGKYA RMRCGAHTNDV KOLITEAVQKI TTESIVIWKG TPKFKLPIQK ETWETTWWTIEY WOATWIPMEW FWNTPLVVLK WYOLEKEP IV GAETFYVDGA 600
 ANRETKLGKA GYVTNRGRQK VVTLTDTNQ KTELQAYILA LDQSGLEVNI VTDSQYALGI IQAQPDQSES ELVNQ.LLEQL IKKEKVYLW VPALKGIGGN 700
 EQVDKLVASAG IRKVIFLDGI DKAQDEHEKY HSNWRAMASD FNLPVVAKA IVASCDKCOL KGEAMHGQVD CSPGIWOLD CTHLEGKVLV AVHVASGYIE 800
 AEVIPAETGQ ETAYFLKLIA GRMPVKTIIH DNGSNFTGAT VRA�WAGI KQEFQGIPYNP QSQGVYVESMN KELKKLIGQV RDQAELHLKTA VQMAVFIHNF 900
 KRKGIGGYS AGERIVDIIA TDIQTKELOK QITKIQNFRV YYRDSRNPLW KGPALKLWKG EGAVTIQDNS DIKVUPRRKA KIIRDYQKQM AGDDCVASRQ 1000
 DED 1003

Pol p10 Protease

PQVTLWQRLP VTIKLGQLK EALLDTGADD TVLEEMSLPG RWKPKMIGGI GGFIKVQRQYD QILIEICGHK AIGTVLVGPV PVNIIGRNLL TQIGCTLNF 99

Pol p66 Reverse Transcriptase (RT/RNase)

PISPIETVPV KLKPQMDGPK VKQWPLTEEK IKALVEICTE MEKEGKISKI GPENPYNTPV FAIKKKDSTK WRKLVDFREL NKRTQDFWEV QLGIPHPAGL 100
 KKKKSVTVLD VGDAYFSVPL DEDFRKYTAF TIPSINNETP GIRYQYNVL P QGWKGSPAIIF QSSMTKILEP FRKQNPDIVI YQYMDDLVVG SDLEIGQHRT 200
 KIEELRQHLL RWGLITTPDKK HQKEPPFLWM GYELHPDKWT VQPIVLPEKD SWTVNDIQLV VPKLNWASQI YPGIKVRQLC KLLRGTKALT EVIPLTEAE 300
 LEIAENREIL KEPVHGVYYD PSKDLIAEIQ KQGQQGWTYQ IYQEPFKNLK TGKYARMRG AHTNDVKQLTE AVQKITTESSI VIWGKTPKFK LP1QKETWET 400
 WWTEYWQATW IPEWEFVNTP PLVKLWYQLE KEP1VGAETF YDGAANRET KLGKAGYVTIN RGRQKVVTLT DTINQKTELQ AIYLALQDSG LEVNIVTDSQ 500
 YALGIIQAZQDQSESELVNLQ IIEQLIKKEK VYLAWPAHK GIGGNEQVDK LVSAGLRKVL 560

Pol p51 RT

PISPIETVPV KLKPQMDGPK VKQWPLTEEK IKALVEICTE MEKEGKISKI GPENPYNTPV FAIKKKDSTK WRKLVDFREL NKRTQDFWEV QLGIPHPAGL 100
 KKKKSVTVLD VGDAYFSVPL DEDFRKYTAF TIPSINNETP GIRYQYNVL P QGWKGSPAIIF QSSMTKILEP FRKQNPDIVI YQYMDDLVVG SDLEIGQHRT 200
 KIEELRQHLL RWGLITTPDKK HQKEPPFLWM GYELHPDKWT VQPIVLPEKD SWTVNDIQLV VPKLNWASQI YPGIKVRQLC KLLRGTKALT EVIPLTEAE 300
 LEIAENREIL KEPVHGVYYD PSKDLIAEIQ KQGQQGWTYQ IYQEPFKNLK TGKYARMRG AHTNDVKQLTE AVQKITTESSI VIWGKTPKFK LP1QKETWET 400
 WWTEYWQATW IPEWEFVNTP PLVKLWYQLE KEP1VGAETF YDGAANRET KLGKAGYVTIN RGRQKVVTLT DTINQKTELQ AIYLALQDSG LEVNIVTDSQ 500
 YALGIIQAZQDQSESELVNLQ IIEQLIKKEK VYLAWPAHK GIGGNEQVDK LVSAGLRKVL 440

Pol p15 RNase

YVDGAANRET KLGKAGYVTIN RGRQKVVTLT DTINQKTELQ AIYLALQDSG LEVNIVTDQV DQSESELVNLQ IIEQLIKKEK VYLAWPAHK 100
 GIGGNEQVDK LVSAGLRKVL 120

Pol p31 Integrase

FLDGIDKAQD EHEKXHSNWR AMASDFNLPP VVAKETIVASC DKCQLGEAM HGQVDCSFGI WOLDOTHLEG KVILVAVHVA SGYIEAEVIP AETGQETAYF 100
 LLKLAGRWPV KTIHTDNGSN FTGATVRAAC WMAGIKQEFG IPYNPQSQV VESMNRKELRK IIQGQRDOAE HLKTAQOMAV FLHNFKRKGG IGGYSAGERI 200
 VDIATDQT KELQKOF-TKI QNFRVYYRDS RNPLWKGPCK LLMKGEGAVV IQDNSDIKVV PRRKAKIIRD YGKQWAGDDC VASRQDED 288

Numbering Positions in HIV

Vif	MENRWQVMIV WQVDRMRIRT WKSLVKHHMY VSGKARGWFY RHHYESPHPR ISSEVHILG DARLVITTYW GLHTGERDWH LGQGVSIIEWR KKRYSTQVDP 100 ELADQLIHLY YFDCCFSDSAI RKALLGHIIS PRCEYQAGHN KVGSLOYLAL AALITPKKKK PPLPSVTKLTT EDRMNKPQKT KGHRGSHTMN GH 192
Vpr	MEQAPEDQGP QREPHEWTL ELLEELKNEA VRHFPRIWHL GLGQHIYETY GDTWAGVEAI HXB2 frameshift \/ MEPVDPRLEP WKHPGSQPKT ACTNCYCKKC CFHQCVCFIT KALGISYGRK KRRQRRAHQ NSQTHQASLS KQPTSOPRGD PTGPKE\$KKK VERETETDPF 101 D
Tat (premature HXB2 stop codon indicated by \$)	HXB2 frameshift \/ Primary splice site MAGRSGDSDE ELIRTVRLIK LLYQSNPPPN PEGTRQARRN RRRWRERQR QHSISERIL GTYLGRSAEP VPLQLPPLER LTLDQNEDCG TSGTQGVGSP 100 QILVESPTVL ESGTKE 116
Rev	Primary splice site Primary splice site TQPIPIVAIV ALVVAIIIAI VVNSIVIIEY RKILRQRKID RLIDRLIERA EDSGNESEGE ISALVEMGVE MGHHAPWDVD DL 82
Vpu (defective start codon)	> gp41 start
Envelope (Env) gp160	Env signal peptide MRVKEKYQHL WRWGWRWGTN LLGMLMICSA TEKLWVTVYY GVPVMKEATT TLFCASDAKA YDTEVHENWYA THACVPTDPM PQEVVLLVNVT ENFNMMWKNDM 100 VEQMHEDTIS LWDQSLKPCV KLTPLCVSLK CTDLKNDNTNT NSSGRMIME KGEIKNCFSN ISTISLRGKVQ KEYAFFYKLD LIPIIDNDTTS YKLTSCTNTSV 200 ITQACPKVSF EPIPHYCAP AGFAALKCENN KTENGTTGPCT NVSTVQCTHG IRPVVSTQLL LNGSIAEEEV VITSVNFNTDN AKTILIVOLNT SVELINCCTRPN 300 NNTRKRTRIQ RGPGRAFVTI GKIGNMRQAH CNISRAKWNN TLKQIASKLQ EQFGNNKTLI FKQSSGGDPE IVTHSFNCGG EFFYCNSTSQL FNSTWFNSTW 400 STEGSNNTEG SDTITLPCRI KQLINNMWQKV GKAMYAPPIS QGIIRCSSNIT GLLTRDGGN SNNESETFRP GGGDMRDNWRL SELYKVVK IEPLGVAPTK 500 AKRRVVQREK RAVGIGALFL GFLGAAGSTM GAASMTLTVQ ARQLLSGIVQ OQNNLIRATE AQQHILQLTY WGIKQLOQARI LAVERYLKQDQ QLLGIWGCSG 600 KLICTTAVPW NASWSNKSLE QIWNHHTWME WDREINNYTS LIHSLIEESQ NQQEKNEQEL LEIDLKWAASLW NWFNNTNWLMW YIKLFLIMVG GLVGLRIVFA 700 VLSIVNRVRQ GYSPISFQTH LPTPRGPDRP EGIEEEGGGER DRDRSIRLVN GSLALIWNDL RSLCLFSYHR LRDLILLIVTR IVELLGRRGW EALKYWWNLL 800 QYWSQELKNS AVSLINATAI AVAEGTDRVI EVVQGACRAI RHPIRRQG LERILL 856
Nef (premature HXB2 stop codon indicated by \$)	MGGMKWSKSSV IGMPTVTERM RRAEPAAADRV GAASRDLKEH GAITSSNTAA TNAACAWLEA QEEEYVGFPV TPOVPLRPMT YKAAAVDLSHF LIKEKGGLEGL 100 IHSQRQRDTL DLMITYHTQGY FPD\$QNYTPG PGVRYPLTFG WCYKLVPVEP DKIEEANKGE NTSUJHPVSL HGMDDPEREV LEWRFDSRLA FHHVARELHP 200 EYFKNC 206

HXB2 Nucleotide Sequence Numbering:

> 5' LTR U3 region start
tggaaaggct aattcaactcc caacgaaagac aagatatacc ttatctgtgg atctaccaca cacaaggcta ctcccctgtat tagcagaact acacaccagg 100
gccaggatc agatccac tgaccctttgg atgggtctac aagcttagtc cagttggcc agagaatgtt gaagaaggcca acaaaggaga gaacaccagg 200
ttgttacacc ctgttagcct gcatggaaatg gatgaccgg agagagaagt gttagatgg aggttgaca gccccttagc atttcatac atggcccgag 300
agttgcattcc ggatgtttc aagaactgtt gacatggac ttgttacaag ggacttccg ctggggactt tccagggggg cgtagggactgg 400

5' LTR U3 region end \/ 5' LTR R repeat start
ggatgtggcga gcccctcagat cctgcataata agcagctgct ttttgcctgt actgggtctc tcttggttaga ccagatctga gcctggagc tctctggcta 500

5' LTR R 5' LTR U5
repeat end \/ region start
actaggaaac ccactgttta aggcttcaata aagcttgcct tgatgtttc aatgtatgt tgccctgtatg tttgtgtactt ctggtaacta gagatccctc 600

5' LTR U5 region end <
agacccttt agtcagttgtt gaaaatctct agcgtggcg cccgaacagg gacctgaaag cggaaaggaa acccaggaggaa ctctctcgac gcaaggactcg 700

> Gag p17 start
gcttgcgttga gcgccgcacgg caagaggccg gggccggcga ctggtagta cgccaaaaat ttgtacttagc ggaggcttaga aggagagaga tgggtgcgag 800
agcgtcagta ttaaggcccc gagaatttggat tcgtatggaa aaaatttcgtt taaggccagg gggaaaggaaa aaataataat taaaacatat agtatggca 900
agcggggac tagaaacgtt cgcgtttaat cctggctgt tagaaacatc aagaggctgt agacaaatac tgggacacgtt acaaccatcc cttcagacag 1000
gatcagaaga acttagatca ttatataataa cgttagcaac cctctattgt gtgcataaa ggtatggat aaaaagacacc aaggaaggtt tagacaagat 1100

Gag p17 end \/ Gag p24 start
agaggaagag caaaaacaaa gtaagaaaaa agcacagcaa gcaaggctg acacaggaca cagccatcg gtcagccaaa attacctat atgcagaac 1200
atccaggggc aaatggtaca tcaggccata tcacctagaa cttaaatgc atggtaaaa gtagttagaa agaaggctt cagccagaa gtgataccca 1300
tgtttcagc attatcggaa ggagccaccc cacaaggatt aaacccatg ctaacacatg tggggaca tcaggcgc atgcaaatgt taaaagagac 1400
catcaatgtgg gaagctgtcgg aatggggatag agtgcattca gtgcattcgg ggcctatgg accaggcccg atgagagaac caagggaaag tgacatagca 1500
ggaaactacta gtacccttca ggaacaaata ggtatggatga caataatcc acctatcc acataatcc cttataaaag atggataatc ctggattaa 1600
ataaaatagt aagaatgtt agcccttacca gcattctggc cataagacaa ggacccatgg aaccctttatg agactatgtt gaccgggttct ataaaactct 1700
agagcccgag caagcttcac aggaggtaaa aaattggatg acagatgttccca aatgttcac ccagattgtt agactatgtt aaaaactttt aaaaactgtt 1800

Gag p24 Capsid end \/ Gag p2 start
ggaccagggg ctadactaga agaaatgtatc acagatgtc agggatgg tagggatggc cttttttggc tgaaggaaatg agccaaagttaa 1900

III-107
DEC 98

Numbering Positions in HIV

Gag p2 end \ / Gag p7 Nucleocapsid start
 caaattcagc taccataatg atgcagagaa gcaattttg gaaccaaaga aagattgtta agtgtttcaa ttgtggcaaa gaaggccaca cagccagaaa 2000

ribosome -1 slip Gag to Gag-Pol

Gag p7 nucleocapsid end \ / Gag p1 start
 Pol start >

Gag p1 end \ / Gag p6 start
 ttgcaggccc ccttagaaaaa aggctgttg gaaatgttga aagaaggac accaaatgaa agattgtact gagagacagg ctaattttt aggaaagatc 2100

> Pol protease start Gag p6 end <

tccccctca gaaggcaggag ccgatagaca aggaactgt tcctttaact tccctcagtttgg caacgaccct tcgtcacaat aaagataggg 2300
 gggcaactaa aggaagctct attagataca ggaggcagat atacagtatt tgtggacata aagctatagg tacagtatta gttagtataa 2400
 tatggagttt tatcaaaggta agacagatgt atcagatact catagaatc 2500

Pol protease end \ / Pol p66 and p51 RT start

cataattgga agaaaatctgt tgactcgat tggttgact ttaatttc ccattagccc tatttagact gtaccgactaa aattaaaggcc aggaatggat 2600
 gggccaaag ttaaacaatg gcccatgaca gaagaaaaaa taaaaggcatt agtagaaatt tggacgaga tggaaaaaggaa agggaaaatt tcaaaaatttg 2700
 ggcctgaaaaa tccatacaat actcccgat ttgcctaaaaa gaaaaagac agtactaaat 99gaaaaattt agtagatttcc agagaactta ataaagaac 2800
 tcaagacttc tgggaagttc aatttaggaat accacatccc gcagggttaa aaaaaaaaaa atcgttaacaat 99gaaaaattt agtagatttcc agagaactta ataaagaac 2900
 gttcccttag atgaagactt caggaagttt acitgcattta ccataccttag tataaaacat gagaacccag 99gatttagata tcagttacaat 99gattttccac 3000
 agggatggaa aggattcacca gcaaatattcc aaagttagcat gacaaaatc tttagccctt ttagaggcttt ttagaggcttt aatcccgacat atagttatct atcaatacat 3100
 ggtatgttq tatgttaggt ctgactttaga aataggcag catagaacaa aaatagaggaa gctgttqaa catgtgttqaa ggtgggact tacacacca 3200
 atcggaaaac acctccatcc tttggatgg gttatgttqaaact ccattctgtat aaatggacag tacaggccat agtgctgcca gaaaaagaca 3300
 gctggactgt caatgacata cagaaggatgt tagggaaattt gaattggca agtcagattt accccaggat taaaaggtaagg caattatgtta aactcccttag 3400
 aggaaccaa gcactaaacag actaataatcc aagtaataatcc gaaacacgaa gaaggcagac 99gaaactgtc tagaactgtc aagaacacgt acatggaggatg 3500
 tattatgacc catcaaaggat tttaatagca gaaatacaga agcaggggca agcccaatgg acatatacaa ttatcaaga gtcattaaa aatctggaaa 3600
 cagggaaaata tgcaagaatg aggggtggcc acactaatga tggaaaacaa ttaacagagg cagtgcacaa aataaccaca gaaaggatag taatatggg 3700
 aaagactct aaatttaaac tggccataca aaagggaaatc tggggacaga gtttggacat gtttggccaa gtcattttttt gtcattttttt 3800

Pol p51 end p66 RT continue \ / Pol p15 RNASE H start

aataccccc ccttagtgtaa attatgttac cagtttagaga aagaacccat agtagggaca gaaacccat gtttggatgg ggcagactaac 3900
 aatttagaaaa agcaggatata gtttactaata gaggaaagaca aaaagtgtc accctaactg acacaacaaa tcagaagact gattttatct 4000
 agctttgcag gatttggat tagaagtaaa catagtaaca gactcacaat atgcattttt gatccatccaa gracaacccg atcaaaatgtga atcagatgtt 4100
 gtcaatcaa taatagagca gtttaataaaaa aaggaaaaaa aatgggggggg aaatgggggg aaatgggggg 4200

Pol RNase H, p66 RT end \ Pol p31 Integrase start

tagtcgtgc tggaaatcagg aaagtactat ttttagatgg aatagaataag gcccagaatg aacatggagaa atatcacatg aattggagag caatggctag 4300
 tgatttaac ctgccacccgt tagtagcaaa agaaataatgta gccagctgtg ataaatgtca gctaaaggaa gaagccatgc atggacaatg agactgtagt 4400
 ccggaaatat ggcaactaga ttgtacacat ttgtacacat ttgttatcc 9gttagcaat ggttagcaat catgttagcca gtggatataat aqaaaggaa gttattccag 4500
 cggaaacagg qcaggaaaca gcatattttc ttgtacacat ttgttatcc 9gtggggggg aatcaaggcg gaatttggaa ttccctacaa tcggggatg tagaactatg gaataaagaa 4600
 tacggttagg gccgcctgtt 9gtggggggg aatcaaggcg atcttaagac agcgtacaa atggcgtat tcatccacaa tttaaaaaga aaagggggggaa 4700
 ttaaagaaaa ttataggaca cggtaaaggat tagacataat agcaacagac atacaacta aagaatrraca aaaacaatatttcc 9400
 ttgggggta cagtccaggg gaaaaggatag tagacataat agcaacagac ttggaaaggcc ttctcttgaa ggtggggg gcgttagtaa tacaagataa tagtgacata 5000

> Vif start

aaagtaggtgc caagaaagaa agcaaagatc attagggatt atggaaaca gatggcgtg gatgtttat agacatctact atgaaaaggcc tcattccaa gatggcgtg tggcaatgg 5100
 tgaaaaaggat tagtaaaaca ccatatgtat gtttcaaggaa aactatattgg ggtctgcata caggaaag agactggcat ttgggtcagg gatctccat 5200
 aagtacat cccactggg gatgtctggg atggcataat aatggcataat aacttagcag accaactaat tcatctgtat tacttgtat gttttcaga ctctgctata 5300
 agaatggggg aaaaaggat tagacacaca agtagaccct gaacttagcag aggacataac aaggtaggat ctctacaataa cttagggacta gcaactaaatc 5400
 agaaaggcct tattaggaca catagtttagc ccttaggtgtg aatatacgcc ttggggtaa aatatacgcc ttggggacta gcaactaaatc 5500

> Vpr start

taacaccaa aaagataaaag ccacctttgc cttagttac gaaactgaca gaggatagat ggaacaaggcc ccagaaggaccc aaggccacaa gagggaggccaa 5600

Vif end <

cacaatgaa ggacacttgc gcttttagag gagcttaaga atgaagctgt tagacatttt ccttaggttt ggctccatgg cttagggccaa catatctatg 5700
 aaactttagg ggtatacttgg gcaggagggg aagccataat aagaatctgc caacaactgc tggttatccca ttttcagaat tgggtgtcga catagcagaa 5800

Tat start > Vpr end <

taggcgttac tcgacacagg agagcaaaatggaggcc tagatccatg actagagccc tggaaaggatc caggaaggatc gcctaaaact gctgttacca 5900

Rev start >

atggcttattg taaaaatgt tgcttttcatg gccaaaggtttgg tttcataaca aaaggctttag gcatctccata tggcaggaa aagcggagac agcgacgaag 6000

Tat, Rev exon end \ /Tat, Rev intron > Vpu start (defective ACG start codon)

agctcatcgt aacagtccaga ctcatcaagg ttctctatca aagcgtatg tagtacatg aacggcaaccc ataccaatag tagcaatgt agcattatgt 6100
 ttagcaataa taatagcaat agttgtgtgg tccatagtaa tcatagaata tagggaaaataa ttaagacaaa gaaaaataga caggtaattt gatagactaa 6200

> Env gp160 start, signal peptide

tagaaagagc aagaagacagt ggcaatgaga aatattcaga ctttgtggaa tggggggccatgctccctt gggatgttga 6300

Numbering Positions in HIV

Vpu end, signal peptide end

III-110
DEC 98

Env gp41, gp160 end < Nef start
 cagatagggt tataagaatg gtacaaggag cttgttagagc tattcgccac atacctagaat gaaataagaca 9999tggaa aggattttgc tataagatgg 8800
 gtggcaagtgt gtcaaaaagt agtgtgttgg gatggctcac tgtaaggaa agaaatgagac gagcttagcc aggaggcagat agggtggag cagcatctcg 8900
 agacctggaa aaacatggag caatccaaag tagcaataca gcagctacca atgctgttg tgcctggcta gaaycacaag aggaggagga 9999ttt 9000

 > 3' LTR U3 region
 ccagtccacac ctcaggtaacc tttaagacca atgactaca accgactgt agatcttagc cacttttaa aaaaaagggg gggactggaa gggctaattc 9100

 Nef end <
 actccccaaag aagacaaggat atccttgatc tgtggatcta ccacaccaa ggctacttcc ctgatttagca gaactacaca ccagggccag gggtcagata 9200
 tccactgacc ttggatgtg gctacaaggat agtaccaggat gaggccagata agatagaaga ggcccaataaa gyagagaaca ccaggtgtt acaccctgtg 9300
 agcctgcatg ggtatggatga cccggagaga gaagtgttag agtggaggt tgacagccgc cttagatttc atcacgtggc ccgagagctg catccggagt 9400

Nef premature stop
 cagatccctgc atataaggcag ctgttttttgc cctgtactgg gtctctctgg ttagaccaga tctgaggctctg ggactaactag ggaaccact 9600

 3' LTR U3 region \ / 3' LTR R repeat
 gcttaagct caataaaaggct tgcccttgagt gcttcaagta gtgtgtggcc gtctgttg tgactctgtg aacttagagat ccctcagacc cttttagtca 9700

 3' LTR U5 end <
 gtgtggaaaatctctaggca 9719

III-111
DEC 98