

1 **Genetically flexible but conserved: a new essential motif in the C-ter domain of**
2 **HIV-1 group M integrases**

3

4 Marine Kanja^a, Pierre Cappy^a, Nicolas Levy^b, Oyndamola Oladosu^b, Sylvie Schmidt^c, Paola
5 Rossolillo^a, Flore Winter^a, Romain Gasser^a, Christiane Moog^c, Marc Ruff^b, Matteo Negroni^{a#}, and
6 Daniela Lener^{a#}

7

8 ^a *Université de Strasbourg, CNRS, Architecture et Réactivité de l'ARN, UPR9002, Strasbourg, France*

9 ^b *Chromatin Stability and DNA mobility, Department of Structural Biology and Genomic, IGBMC,*
10 *Strasbourg University, CNRS, INSERM, Illkirch, France*

11 ^c *Molecular Immuno-Rheumatology Laboratory, UMR1109, FMTS, Université de Strasbourg, INSERM,*
12 *Institut de Virologie, Strasbourg, France*

13 Running head: HIV integrase new functional motif CTD

14

15 # Corresponding authors: Daniela Lener and Matteo Negroni

16 E-mailing addresses: d.lener@ibmc-cnrs.unistra.fr; m.negroni@ibmc-cnrs.unistra.fr

17 Mailing address: Institut de Biologie Moléculaire et Cellulaire

18 2 allée Konrad Roentgen, 67084 Strasbourg Cedex, France

19

20 Abstract word count number: 219 (importance 116)

21 Main text word count number: 10 056

22 **ABSTRACT**

23 Using coevolution-network interference based on the comparison of two phylogenetically distantly
24 related isolates, one from the main group M and the other from the minor group O of HIV-1, we
25 identify, in the C-terminal domain (CTD) of integrase, a new functional motif constituted by four non-
26 contiguous amino acids (N₂₂₂K₂₄₀N₂₅₄K₂₇₃). Mutating the lysines abolishes integration through
27 decreased 3'-processing and inefficient nuclear import of reverse transcribed genomes. Solution of the
28 crystal structures of wt and mutated CTDs shows that the motif generates a positive surface potential
29 that is important for integration. The number of charges in the motif appears more crucial than their
30 position within the motif. Indeed, the positions of the K could be permuted or additional K could be
31 inserted in the motif, generally without affecting integration *per se*. Despite this potential genetic
32 flexibility, the NKNK arrangement is strictly conserved in natural sequences, indicative of an effective
33 purifying selection exerted at steps other than integration. Accordingly, reverse transcription was
34 reduced even in the mutants that retained wt integration levels, indicating that specifically the wt
35 sequence is optimal for carrying out the multiple functions integrase exerts. We propose that the
36 existence of several amino acids arrangements within the motif, with comparable efficiencies of
37 integration *per se*, might have constituted an asset for the acquisition of additional functions during
38 viral evolution.

39 **IMPORTANCE** Intensive studies on HIV-1 have revealed its extraordinary ability to adapt to
40 environmental and immunological challenges, an ability that is also at the basis of antiviral treatments
41 escape. Here, by deconvoluting the different roles of the viral integrase in the various steps of the
42 infectious cycle, we report how the existence of alternative equally efficient structural arrangements for
43 carrying out one function opens on the possibility of adapting to the optimisation of further
44 functionalities exerted by the same protein. Such property provides an asset to increase the efficiency
45 of the infectious process. On the other hand, though, the identification of this new motif provides a
46 potential target for interfering simultaneously with multiple functions of the protein.

47

48 Introduction

49 Integration of reverse transcribed viral genomes into the genome of the infected cell is a peculiar
50 feature of the replication strategy of retroviruses, carried out by the viral enzyme integrase (IN) in a
51 two-step reaction. In HIV-1, after the achievement of DNA synthesis in the cytoplasm of the infected
52 cell, it first catalyses the removal of a conserved GT dinucleotide from the 3' ends of the viral DNA (3'
53 processing), leaving CA_{-OH} 3' ends bound to the active site. Subsequently, once the viral DNA has
54 been imported in the nucleus, the reactive CA_{-OH}-3' ends attack the cellular DNA leading to the
55 generation of the provirus (1, 2).

56 Besides this enzymatic function, HIV-1 IN is involved, through non catalytic activities, in several other
57 steps of the viral replication cycle. As a component of the Gag-Pol polyprotein precursor, it participates
58 in Gag-Pol dimerization, essential for the auto-activation of the viral protease and, consequently, for
59 viral particle maturation (3-5). During capsid morphogenesis, it is involved in the recruitment of the
60 genomic RNA inside the core of the viral particle (6). As a mature protein, it interacts with the viral
61 polymerase (reverse transcriptase, RT) to optimize reverse transcription of the viral genome (7-9).
62 Through the interaction with the cellular protein LEDGF/p75, it targets actively transcribed genes as
63 sites for integration (10). Finally, as a component of the pre-integration complex (PIC), IN is also
64 involved in nuclear import of the reverse transcription product, a peculiar feature of lentiviruses that
65 allows the infection of non-dividing cells.

66 This ability relies on the virus capacity to enter the nucleus *via* an active passage through the nuclear
67 pore complex (NPC) (11, 12). Several lines of evidence have indicated that the capsid (CA) protein is
68 crucial for nuclear entry (13, 14), through its interaction with several nucleoporins (Nups) forming the
69 NPC (Nup 358, Nup 153, Nup 98) (15-17) and with the transportin-3 (18, 19). Nevertheless, several
70 studies have indicated that IN has karyophilic properties. Namely, it contains a basic bipartite nuclear
71 localization signal (NLS) (20) as well as an atypical NLS (21), and also binds several cellular nuclear
72 import factors. Interactions of IN with importin α/β (22), importin 7 (23), importin $\alpha 3$ (24), Nup153 (25),
73 Nup62 (26) and transportin-3 (27-29) have been documented. Indeed, the mutation of amino acids,
74 mostly located in the C-terminal domain of the IN, responsible for binding to nuclear import factors,
75 results in non-infectious viruses impaired in nuclear import (23, 24, 28, 30).

76 The functional form of the HIV-1 integrase is made up of a dimer of dimers, which assemble in highly
77 ordered multimers of these tetramers (31, 32). Three domains, connected by flexible linkers, constitute
78 HIV-1 IN: the N-terminal domain (NTD), the catalytic core domain (CCD) and the C-terminal domain
79 (CTD) (33, 34). While the NTD is mostly involved in protein multimerization (35, 36), the CCD is mostly
80 responsible for catalysis, and for binding to the viral and cellular DNA as well as to the cellular cofactor
81 LEDGF (37-40). Finally, the CTD is involved in DNA binding during integration (41, 42), in protein
82 multimerization (35), in the interaction with the reverse transcriptase (7, 9) and in the recruitment of
83 the viral genomic RNA (gRNA) in the viral core (6). Overall, the intrinsic flexibility of the protein, the
84 multiple steps required to achieve integration, and the multimeric nature of the integration complex
85 make the involvement of the different parts of the protein in the various functions of the integrase very
86 complex and still not fully elucidated.

87 In addition, the multiple tasks that the IN must accomplish during the infectious cycle and the
88 complexity of its supramolecular structures are expected to impose functional constraints that
89 ultimately may limit its genetic diversity. Retention of functionality despite sequence variation strongly
90 relies on covariation, inside or outside the mutated protein. When an initial mutation negatively alters
91 the protein functionality, compensatory mutations can restore it, at least partially. Therefore, the
92 sequences of homologous proteins in different HIV variants are the result of independent evolution
93 pathways, with independent covariation networks specifically generated for each pathway. Chimeric
94 genes between variants of a given protein can perturb such networks and result in the production of
95 non-functional proteins. This information can then be exploited to probe the existence of functional
96 motifs in proteins. For considerably divergent viruses, as those derived from independent zoonotic
97 transmissions, this approach can be particularly powerful. This is the case for HIV-1 groups M and O
98 that derive from simian viruses infecting chimpanzees and gorillas, respectively. Here, we exploit the
99 natural genetic diversity existing between these groups to generate chimeric integrases. A detailed
100 characterization of the individual amino acids that differ in the non-functional chimeras has then led to
101 the identification and functional characterization of a new motif, in the CTD of HIV-1 group M
102 integrase, essential for viral integration.

103

104

105 **Results**

106 ***Analysis of intergroup M/O chimeras in the CTD of IN***

107 The functionality of the integrases studied in this work was evaluated following the protocol outlined in
108 Figure 1A-B and detailed in Materials and Methods. For this, we replaced the original RT and IN
109 sequences of the p8.91-MB (see Materials and Methods) by those of either one isolate of HIV-1 group
110 M, subtype A2, referred herein as "isolate A" or one isolate of HIV-1 group O, referred herein as
111 "isolate O". The resulting vectors were named RTA-INA (vRTA-INA) and (vRTO-INO), respectively.
112 With these vectors, we estimated the functionality of the integrases by measuring the efficiency of
113 generation of proviral DNAs. Since the number of proviral DNAs generated for each sample is
114 dependent not only on the levels of functionality of the integrase but also on the amount of total viral
115 DNA generated after reverse transcription, we estimated the amount of total viral DNA generated by
116 each sample by qPCR as described in Materials and Methods. In parallel, we measured the amount of
117 proviral DNA generated either by the puromycin assay or by the Alu qPCR assay as described in
118 Materials and Methods (Evaluation of integration by puromycin assay). The amount of proviral DNA
119 divided by that of total viral DNA provides an estimate of the efficiency of integration. Comparable
120 efficiencies of integration were measured with the two vectors, irrespective of whether the estimation
121 was done using the puromycin assay ($71 \pm 13 \%$ and $72 \pm 24 \%$ the level of the reference vector
122 v8.91-MB respectively) or the Alu qPCR assay ($70 \pm 14 \%$ and $68 \pm 15\%$, respectively, Figure 1C).
123 Throughout this work, the efficiency of integration has always been evaluated by the puromycin-
124 resistance assay normalised by the amount of total DNA. Control vectors, in which the catalytic activity
125 either of the integrase or of the reverse transcriptase have been abolished, in vRTA-INA, by the
126 introduction of the D116A mutation in IN or of the D110N-D185N mutations in RT (43, 44), gave the
127 expected results (Figure 1C).

128 We chose to probe the existence of functional motifs in the C-terminal domain of integrase because
129 this domain is involved in several non-catalytic functions of the protein. The C-terminal domain of INO
130 used in this study is 10 amino acids longer (212-298) than that of INA (212-288, Figure 2A). We

131 constructed three chimeras between isolates A and O, named after the position, in amino acids from
132 the beginning of the IN-coding region, where the sequence shifts from that of one isolate to that of the
133 other (Figure 2B). Chimera A(1-212)-O(213-298) is constituted by INA with the entire CTD from INO;
134 chimera A(1-285)-O(286-298) is INA with the additional 10 amino acids of INO at the C-ter end plus
135 the two most C-ter different amino acids; finally, as the region between position 212 and 288 differs in
136 12 amino acids, chimera A(1-272)-O(273-298) was constructed in such a way as to split the 12
137 different amino acids in two groups of 6.

138 We first performed western blots (Figure 2C) on viral particles to monitor the degree of proteolytic
139 processing of the Gag precursor (Pr55Gag), since incomplete processing would result in immature
140 viral particles, affecting infectivity. No significant differences in Pr55Gag were observed between
141 isolate O and chimerical constructs compared to isolate A (Figure 2D). We then evaluated the
142 efficiency of reverse transcription (measuring the amount of viral DNA produced by qPCR) and of
143 integration (as described above). Only chimera A(1-272)-O(273-298) exhibited significant defects in
144 both reverse transcription and integration (Figure 2E-F), suggesting that a covariation network,
145 present between positions 212 and 285, was broken in this chimera. Since in these experiments the
146 IN is expressed from p8.91-MB and not from the genomic RNA, it can be ruled out that the
147 phenotypes observed are due to an effect of the mutations on the genomic RNA, as for example on
148 the process of splicing, as it has been previously described for some mutants of the C-terminal domain
149 of IN (45).

150

151 ***Characterization of IN CTD***

152 In order to evaluate the individual contribution of the 10 amino acids differing between positions 212
153 and 285 (Figure 2A), each residue in IN A was individually replaced by those of IN O and the ten-point
154 mutants were tested for processing of Pr55Gag, reverse transcription and generation of integrated
155 proviruses.

156 Except for mutant N254K, no significant difference in level of Pr55Gag proteolytic processing was
157 observed between mutants and parental vector A (Figure 3A). The effect on reverse transcription was

158 an overall reduction of efficiency for most of the mutants, with a residual efficiency between 45 and
159 90% that of the parental vector A (Figure 3B). Concerning integration efficiency, instead, the majority
160 of the mutants did not show a significant decrease, except for mutants K240Q and K273Q for which
161 integration was dramatically impaired (Figure 3C). This suggests a specific implication of these two
162 residues in the integration process. When the two mutations were combined (K240Q/K273Q mutant),
163 while the level of reverse transcription remained above 40 % that of wt IN A, integration dropped to
164 undetectable levels (Figure 3D).

165 To discriminate between the role of the charge of K₂₄₀ and K₂₇₃ from that of their possible acetylation,
166 we replaced both residues by two R (K240R/K273R mutant). The level of integration of this mutant
167 was comparable to that of wt IN A (Figure 3E), indicating that the presence of a positive charge and
168 not acetylation at these positions was important for integration efficiency. However, these mutations
169 reduced by half both Pr55Gag proteolytic processing and reverse transcription (Figure S1A).

170 In both mutants showing a marked defect in integration (K240Q and K273Q), a K (positively charged
171 polar side chain) was replaced by a Q (non-charged polar side chain), the amino acid present in
172 isolate O at the corresponding positions. Conversely, in isolate O, two K are present in positions
173 where a polar non charged amino acid (N in both cases) is present in isolate A (positions 222 and 254,
174 Figure 2A). Therefore, in order to evaluate if also the two non-charged polar amino acids (N) present
175 in isolate A at positions 222 and 254 are essential, we replaced them by a non-polar amino acid like
176 leucine (mutant LKLK) and, in parallel, by a non-charged polar residue, Q (mutant QKQK). While in
177 the LKLK mutant the efficiency of integration dropped to almost undetectable levels, in the QKQK one
178 it was comparable to that of the wt enzyme, suggesting that the presence of a polar residue at these
179 positions is essential (Figure 4A). To understand whether the polar nature of the amino acid at
180 positions 222 and 254 is enough to retain functionality, the N were replaced by two threonine, which
181 are polar but do not have the amide group of asparagine. In this case (TKTK mutant) integration
182 dropped to undetectable levels (Figure 4A) indicating that not only the polarity is important but also the
183 functional group carried by the amino acid. Therefore, the biochemical features of all four residues
184 identified are important.

185 Finally, we wondered whether the residues present at positions 222, 240, 254 and 273 could be
186 interchanged between isolates O and A. Therefore, we generated the quadruple mutant of isolate A
187 N222K/K240Q/N254K/K273Q (called KQKQ for simplicity). Remarkably, the integration efficiency of
188 this mutant was not significantly different from that of wt IN A (Figure 4B), indicating the existence of a
189 functional link between these four positions.

190 The alignment of HIV-1 IN sequences reveals a strong conservation of the amino acids
191 N₂₂₂K₂₄₀N₂₅₄K₂₇₃ in group M (Figure 4C). To confirm the need for K₂₄₀ and K₂₇₃, observed in isolate A,
192 also for other isolates of group M, we introduced the K240Q/K273Q double mutation (NQNQ mutant)
193 in integrases from three other primary isolates of group M (Figure 4D). In all cases, a dramatic drop in
194 integration was observed with respect to the corresponding wt integrases, confirming the results
195 obtained with isolate A. The importance of the two K in the motif was therefore confirmed in isolates
196 from the most widespread HIV-1 group M subtypes in the epidemics, subtypes A, B, C, and CRF02
197 being responsible for 79% of the HIV-1 infections worldwide (46).

198 The possibility of permuting the positions of the four amino acids at positions 222, 240, 254 and 273
199 indicates a functional relationship between these residues that can therefore be considered as a
200 functional motif that, based on the identity of the amino acids present at these positions in isolates of
201 group M, we refer to as the "NKNK" motif.

202

203 ***Importance of the lysines in the NKNK motif***

204 To understand to which extent the number and the positions of the K in the motif influence IN
205 functionality, we generated a series of mutants based on the replacement of the amino acids present
206 in isolate A by those of isolate O. We thus tested all possible variants (Figure 5A) containing either
207 only one K (four mutants, Figure 5B), two K (five mutants plus the wt, Figure 5C), three K (four
208 mutants, Figure 5D) or four K (one mutant, KKKK, Figure 5A) at any of the positions in the motif. The
209 presence of a single K led on average to a drop to 20% of integration with respect to wt IN A, whereas
210 when two or more lysines were present in the motif, levels of integration were close to those of wt IN
211 A, ranging from 75 to 137% (Figure 5A). The mutant with no K, where the motif sequence has been

212 changed from NKNK to NQNQ, confirmed the total loss of integration already observed with this
213 mutant (see Figure 3D). Finally, from the analyses of the different mutants it appears that the
214 presence of a K at the first position of the motif (position 222) consistently leads to a higher level of
215 integration in all classes of mutants (those with 1, 2, or 3 K). Interestingly, though, position 222 has a
216 N in the wt enzyme.

217 When considering individual mutants within the different classes, we observed a significant decrease
218 in functionality for all the mutants possessing only one K (Figure 5B). For the mutants containing two
219 K, three variants were at least as functional as wt IN A (NKNK in the figure), while two displayed a
220 significant reduction (Figure 5C). Finally, all mutants containing three K were at least as functional as
221 the parental IN A (Figure 5D). Remarkably, the results obtained with the mutants containing three or
222 four K indicate that the positively charged residues can replace the polar ones, while the reverse is not
223 the case, as shown by the mutants with none or only one K. Overall, these results indicate that at least
224 two K are required to have wt levels of integration, even if not all the positions in the motif are
225 equivalent. Instead, all mutants impacted reverse transcription with a reduction to 40-80% of the wt IN
226 A (Figure S2).

227

228 ***The NKNK motif in replication-competent viruses***

229 To confirm the observations made in the single infection cycle system, some mutants were then tested
230 in a replication-competent system using NL4.3 as primary virus. Mutants of the class containing two K
231 in the motif (the number of K found in circulating viruses) and with a marked phenotype were chosen
232 for this analysis. Besides the wt A sequence, we chose three mutants that either retained integration
233 (KQKQ and KQNK) or exhibited reduced integration (NQKK) (Figure 5C). To construct the four
234 variants, we replaced the sequence of NL4.3 CTD by that of isolate A, either wt or carrying the KQKQ,
235 KQNK or NQKK motifs (Figure 5E).

236 The infectivity of the virus carrying the whole CTD of INA instead of that of NL4.3 (called NL4.3 CTD
237 A, Figure 5E) was comparable to that of wt NL4.3 virus, set as reference, indicating that the
238 replacement of the whole CTD from NL4.3 by that of isolate A did not impact viral infectivity (Figure

239 5F). Regarding the mutants, the results well recaptured the observations made with a single infection
240 cycle (Figure 5C): the infectivity was maintained for KQKQ and KQNK mutants while it was markedly
241 decreased with NQKK motif (Figure 5F).

242

243

244 ***Role of the lysines of the motif in the integration process***

245 In order to characterise in which steps of the infectious cycle are involved the lysines of the NKNK
246 motif, we evaluated the effect of their mutation in two steps (other than reverse transcription) upstream
247 the integration of the pre-proviral DNA in the chromosomes of the host cell. In particular, by
248 quantifying the two LTR circles (2LTRc), we evaluated nuclear import and, by characterizing the LTR-
249 LTR junctions of 2LTRc, the efficiency of 3' processing, which takes place in the cytoplasm, before
250 nuclear import.

251 2LTRc are exclusively formed in the nucleus and are, therefore, useful markers for nuclear import of
252 the reverse transcribed genomes (47). They are generated when the full-length reverse transcription
253 products are not used as substrate for integration. If a mutant is defective in catalysis but carries out
254 nuclear import efficiently (as mutant D116A), 2LTRc should accumulate with respect to a wt IN.
255 Instead, if the mutant is also impaired in nuclear import, 2LTRc will either not increase with respect to
256 the wt IN or increase but more modestly than for D116A.

257 Hence, to monitor nuclear import, we measured the amount of (2LTRc) in wt IN A, in mutants
258 containing either no K (NQNQ) or only one (either K₂₇₃, NQNK mutant, or K₂₄₀, NKNQ mutant). IN A
259 D116A mutant was used as a control. This mutant being totally inactive for integration was considered
260 to produce the highest accumulation of 2LTRc, set at 100%. As expected, the level of 2LTRc found
261 with wt IN A, which efficiently imports and integrates the reverse transcribed genome, was significantly
262 lower (25%) than that of the D116A mutant. As shown in Figure 6A, despite their inability to generate
263 proviral DNA, the mutants had levels of 2LTRc significantly lower than IN A D116A, indicative of a
264 defect in nuclear import.

265 To estimate the efficiency of nuclear import in the mutants, we estimated the level of 2LTRc (Table 1,
266 line 2, "theoretical level"), that could be obtained if no defect in nuclear import was present. We then
267 calculated the efficiency of nuclear import as the ratio between the level of 2LTRc observed
268 experimentally (Table 1, line 3) and the theoretical one. If no defects in nuclear import are present,
269 ratios should be around 1, while defects in nuclear import would yield ratios <1 . The ratios found for
270 the three mutants were in the 0.31-0.35 range (Table 1, line 4), indicative of a reduction of nuclear
271 import to approximately 1/3 that of the wt enzyme. Therefore, the defects in nuclear import contribute
272 to the decrease in integration found with these mutants, but cannot alone account for the low levels
273 observed, particularly in NQNQ and NQNK mutants for which integration was almost undetectable
274 (Table 1, line 1).

275 The efficiency of 3' processing carried out by IN was then analysed by quantifying the different LTR-
276 LTR junctions in the 2LTRc. The 2LTRc found in the nucleus are generated from DNAs carrying either
277 unprocessed or processed 3' ends. In the first case, the 2LTRc will present "perfect junctions" (PJ)
278 while in the second the junctions will be "imperfect". A high ratio of PJ/2LTRc is therefore indicative of
279 inefficient 3' processing. We found that mutating both K (mutant NQNQ) or only K₂₄₀ (mutant NQNK)
280 led to results not significantly different from those obtained with the IN A D116A catalytic mutant
281 (Figure 6B), indicative of a marked defect in 3' processing. Mutating K₂₇₃ (mutant NKNQ), instead, did
282 not affect the process, with PJ/2LTRc values comparable to those of the wt enzyme.

283 To evaluate the contribution of defects in 3' processing to the decreased integration efficiency
284 observed with the various mutants, we first estimated the maximum diminution of PJ/2LTRc ratio
285 observed for a fully competent enzyme (wt IN A). The PJ/2LTRc ratio for wt IN A was 0.54 that of IN A
286 D116A (Table 1, line 5), corresponding to a reduction of 46 % due to 3' processing (Table 1, line 6).
287 The ratio PJ/2LTRc with respect to IN A D116A was then calculated for each mutant and the resulting
288 value was divided by 0.46, obtaining an estimate of the efficiency of 3' processing relative to that
289 observed for wt IN A (Table 1, line 7). 3' processing of NQNQ and NQNK mutants was dramatically
290 reduced, 15 % and 28 % that of wt IN A, respectively. Mutating K₂₇₃, instead, only decreased 3'
291 processing to 74 % that of wt IN A.

292 Finally, in order to understand if these two types of defects (nuclear import and 3' processing) were
293 sufficient to explain the integration defects observed with the various mutants, we combined the effect
294 of these defects (Table 1, line 8). The values obtained account remarkably well for the efficiencies of
295 integration observed (lines 1 and 8 of Table 1) indicating that the decrease observed when mutating
296 the K of the NKNK motif, once normalized for the differences observed in the amount of viral DNA
297 produced, is essentially due to alterations in these two processes.

298

299

300

301 ***Structural analysis of wt and mutant integrase C-ter domains***

302 To understand the structural bases for the functional differences observed in the NKNK motif mutants,
303 the crystal structures of the C-terminal domain (IN CTD, 220-270) of wild type IN A and of the
304 reference strain NL4.3 were solved at 2.2 Å and 1.3 Å of resolution, respectively. For both structures,
305 K₂₇₃ was not included as it is in a disordered region of IN. For all crystal forms we observed a strong
306 packing interaction through the His-tag coordinating a Nickel ion (Figure S3A). The quality of the
307 structures and maps is shown in Figure S3, panels B-D. The structures had the same topology,
308 consisting in a five-stranded β-barrel (Figure S4A). The region encompassing the positions of the motif
309 (Figure 7A-B) generates a surface endowed with a positive potential (circled in yellow in Figure 7C-D),
310 suggesting that this feature could be important for the functionality of the IN. In this case, it is expected
311 that inserting additional lysines in the motif (as for the mutants containing three or four lysines) would
312 retain functionality and, conversely, removing the K (mutants with one K or no K) would affect it. This
313 is what we observed in Figure 5. Nevertheless, the correlation between surface potential and
314 functionality is less clear for the mutants where the number of K in the motif (two) is not altered but
315 their positions are permutated with polar amino acids.

316 To clarify this point, we solved, at a 2.0 Å resolution the crystal structure of the CTD of the NQKK
317 mutant, which was the one displaying the most dramatic drop in integration among the mutants
318 possessing two K (Figure 5C). The NQKK CTD crystalized in a different space group and had three

319 chains in the asymmetric unit. The superposition of the five structures (Figure S4B) corresponding to
320 NL4.3 CTD, to A CTD and to the three molecules in the asymmetric unit for the NQKK CTD (chains A,
321 B and C) did not show significant differences in the main chain fold (Root-mean-square deviation of
322 atomic positions, RMSD: NL4.3 CTD vs A CTD = 0.395 Å, A CTD vs NQKK CTD ABC = 0.653 Å,
323 NQKK CTD chain A vs chain B vs chain C = 0.558 Å). Interestingly, the positive surface electrostatic
324 potential observed for the wt enzyme was markedly perturbed in the NQKK mutant (yellow circle in
325 Figure 7E-F), a change that could well account for the decrease of functionality of the NQKK mutant
326 integrase.

327 To further analyse the impact of the mutations on the structure, we defined the regions which are
328 naturally disordered (Intrinsically Disordered Regions, IDRs) by the superposition of the three
329 molecules in the asymmetric unit of the NQKK CTD structure. We assume that the change in the
330 RMSD obtained among the three molecules represents the natural IDRs. For the main chain,
331 disordered regions with high RMSD are 228-232 and 243-248 (Figure S5A). These same regions are
332 found to be similarly disordered when comparing the main chain RMSD of the C-terminal domain of IN
333 A to that of IN A NQKK chains A, B and C (Figure S5B), indicating that the mutations have no effect
334 on the C-alpha backbone fold of IN CTD.

335 The three structures were then analysed from the standpoint of the arrangement of the side chains.
336 Calculating side chains RMSD, disordered portions were found to correspond to regions 222-225,
337 228-232 and 243-248 (Figure S5C). The comparison of A CTD and NL4.3 CTD, which differ between
338 positions 220-270 only by a single amino acid change (V234 in NL4.3 replaced by I234 in A),
339 expectedly, did not reveal significant differences between the two structures. Instead, when comparing
340 A CTD to NQKK CTD A, B and C side chain deviation, we observed a difference in the structure of
341 region 235-237 (Figure S5D). This difference appears to be due to the N254K mutation that induces a
342 displacement of the side chain of lysine 236 (white arrows in Figure 7G). This is likely due to a
343 repulsive interaction between the side chains of the two lysines, resulting in a perturbation of the
344 structure in the 235-237 region.

345 To evaluate the importance of charge configurations in the context of the C-terminal domain of
346 retroviral integrases, we performed an analysis of the electrostatic charge surface potential for other

347 lentiviruses as well as for other retroviruses. Integrases C-terminal domain structures are available for
348 HIV-1 A2, PDB 6T6I (this publication); HIV-1 PNL4.3, PDB 6T6E (this publication); SIV, PDB 1C6V
349 (48); MVV, PDB 5LLJ (49); RSV, 1C0M (50); MMTV, PDB 5D7U (51); MMLV, PDB 2M9U (52); PFV,
350 PDB 4E7I (53). First, we superposed the available structures. The superposition shows that they
351 share a common fold (Figure 8A) as well as a low Root Mean Square Deviation on secondary
352 structure backbone despite a very low sequence identity for some cases (Table S1). A structure-based
353 sequence alignment has then been performed (Figure 8B). Surprisingly, despite a low overall
354 sequence identity (10 – 20 %) for some integrases, several regions have a strong local sequence
355 similarity (red and yellow background in Figure 8B) while no conservation is observed at positions 222,
356 240 and 254 among lentiviruses nor retroviruses. However, when we compared the electrostatic
357 surface potentials of all structures (Figure 8, Panels C-L), we could define two general retroviral
358 classes. A first class represented by lentiviruses where the surface corresponding to the one delimited
359 by the NKNK motif in HIV-1 M is basic and a second class represented by the other retroviruses
360 tested (orthoretroviruses α , β , γ , and spumaretroviruses) where this surface is acidic or neutral.

361

362

363 **Discussion**

364 Here, by performing a systematic comparison between the non-conserved amino acids in the CTD of
365 the HIV-1 group M and group O integrases, we identify, in group M, a highly conserved motif that is
366 essential for integration. The motif is constituted of two asparagines and two lysines (N₂₂₂K₂₄₀N₂₅₄K₂₇₃)
367 all required for the generation of proviral DNA. In particular, when the K were mutated, integration was
368 abolished due to the cumulative effects of decreased 3' processing and nuclear import of the reverse
369 transcription products (Table 1). Replacing the K by R did not affect integration (Figure 3E),
370 suggesting that the essential feature of the K is their positive charge. Importantly, the positions of the
371 two K of the motif could be permuted without affecting the functionality of the integrase in most
372 cases (Figure 5).

373 A potential explanation for the retention of functionality when permutating the positions of the K across
374 the motif comes from the structural data on the CTD obtained in this work. We have solved the crystal
375 structure of the CTD of the wt IN A used in this study as well as that of the K240Q-N254K mutant
376 (referred as NQKK in the result section). In the structure of the wt enzyme, the residues constituting
377 the motif (except K₂₇₃ which is part of an unresolved region) generate a positively charged surface
378 (Figure 7A-D). This positive electrostatic potential surface is absent, instead, in the NQKK mutant
379 (Figure 7E-F) which, despite the presence of two K, displays a drastic reduction of integration
380 efficiency. These results, combined to the tests of functionality of the different mutants, suggest that
381 the relevant parameter is the presence of a positive charge across this surface. Charged residues
382 have a strong effect on the surface potential. The nature of the amino acid side chain (charged, polar,
383 non-polar) on the surface of the protein defines the surface potential. Charged and polar groups,
384 through forming ion pairs, hydrogen bonds, and other electrostatic interactions, impart important
385 properties to proteins. Modulation of the charges on the amino acids, e.g. by pH and by post-
386 translational modifications, have significant effects on protein – protein and protein – nucleic acid
387 interactions (54). In addition to residues carrying net charges, also polar residues have significant
388 partial charges and can form hydrogen bonds and other specific electrostatic interactions among
389 themselves and with charged residues (55). In the case of the present study, the possible contribution
390 of these mechanisms to the functionality of the integrase could be reflected by the loss of functionality
391 observed by replacing the N, which carries an amide side chain (–CONH₂), by either a non-polar
392 amino acid (L) or by a polar one (T) but carrying a hydroxyl side chain (–OH). The analysis of the
393 electrostatic surface potential for the integrases C-terminal domain of the retroviruses for which this is
394 known showed that, despite a low sequence identity among some of the retroviruses (Table S1), the
395 topology of the structure is maintained (Figure 8A) and the analysis of the surface electrostatic
396 potential splits the viruses studied in two classes. One, constituted by lentiviruses, for which the
397 surface delimited by the NKNK motif of HIV-1 M contains basic charges (in some cases brought by
398 amino acids non corresponding to those of the NKNK motif of HIV-1 M). The second class including
399 the other orthoretroviruses studied (α , β , γ) and spumaretroviruses where this surface contains acidic
400 and neutral regions. This presence of basic regions, specifically in lentiviruses, could contribute to
401 some specificity of lentiviral biology as to increase the efficiency of infection of quiescent cells.

402 The importance for protein functionality of charge configurations and clusters in their three-
403 dimensional structures has been underlined by several studies (56-59). Charge permutations have
404 been used in the NC region of the Gag protein for the Mason-Pfizer Monkey Virus (60). This basic
405 region could be replaced with nonspecific sequences containing basic amino acid residues, without
406 altering its functionality while mutants with neutral or negatively charged residues showed a large drop
407 in viral infectivity in single round experiments. Moreover, a mutant exhibiting an increased net charge
408 of the basic region, was 30% more infectious than the wild type. Also, in our study, increasing the
409 positive charge of the NKNK motif of HIV-1 IN by introducing a third K leads to a slight increase in
410 integration with respect to wt IN (Figure 5A and D).

411 As retention of IN functionality relies on the electrostatic surface potential rather than on the specific
412 positions of the positively charged amino acid, we infer that this region is probably involved in the
413 interaction with a partner carrying a repetitive negatively charged biochemical motif, as the
414 phosphates of the nucleic acids backbone. Alternatively, the partner could be a disordered region of a
415 protein that can rearrange to preserve the interaction when the positions of the positive charges are
416 permuted across the surface of the NKNK motif. Indeed, the molecular recognition between charged
417 surfaces and flexible macromolecules like DNA, RNA and intrinsically disordered protein regions has
418 been observed for the Prototype Foamy Virus and Rous Sarcoma Virus Gag precursors (61, 62), for
419 EBNA proteins of the Epstein-Barr virus (63), UL34 protein of the Herpes Simplex Virus (64) as well as
420 for cellular proteins like APOBEC3G (65). Moreover, the presence of asparagines in the motif, which
421 we show are required for integrase functionality, could contribute to the interaction with the nucleic
422 acid or with a protein partner through hydrogen bonds with the bases (54) or with polar amino acids
423 (55, 66, 67), respectively. The analysis of the motif in the context of the well-characterised structure of
424 the intasome (31), mimicking a post-integration desoxyribonucleic complex, indicates that the residues
425 forming the electrostatic surface point toward the solvent, at the exterior of the structure (Figure S6).
426 This is coherent with the observation that mutating the motif does not affect late steps of the
427 integration process, but rather earlier ones as 3' processing and nuclear import.

428 We show that the NKNK motif of the CTD is involved in 3' processing and nuclear import of the
429 reverse transcription product. Indeed, the removal of the K impacts both processes and when
430 combined, these effects are sufficient to account for the drop of infectivity to the undetectable levels

431 observed in the absence of K in the motif (Table 1). Since mutating K₂₄₀ has a strong impact on both 3'
432 processing and nuclear import, while K₂₇₃ appears to be predominantly involved in nuclear import, it is
433 possible that the involvement of the motif in these two processes implicates structurally distinct
434 functional complexes. This is the first finding of an implication of the HIV-1 IN CTD in 3' processing. So
435 far, only the involvement of the catalytic domain had been demonstrated (39, 68, 69). Since it has
436 been shown that different oligomerization states of IN influence specifically the ability to carry out 3'
437 processing or strand transfer (70), it is possible that the electrostatic surface formed by the NKNK
438 motif help stabilize the oligomeric state that allows 3' processing.

439 Concerning nuclear import, it is known that HIV-1 IN binds several cellular nuclear import factors
440 through basic amino acids of the CTD, and that abolishing these interactions leads to non-infectious
441 viruses displaying a severe defect in nuclear import. Here, we extend the regions of IN involved in this
442 process by describing the need for a new motif, although it cannot be discriminated whether its
443 involvement is direct or mediated by the interaction with a partner protein with karyophilic properties.

444 Some of the residues constituting the N₂₂₂K₂₄₀N₂₅₄K₂₇₃ motif have been previously characterized
445 showing their implication in different steps of the infectious cycle. One is the involvement in reverse
446 transcription. The integrase CTD interacts with the reverse transcriptase to improve DNA synthesis
447 (71, 72). In one study, the double mutation K240A/K244E caused a decrease in reverse transcription
448 to around 20 % the levels of the wt enzyme (72) while the K244E mutation alone caused a reduction
449 of 40 % of RT efficiency (73), suggesting that K₂₄₀ also contributes to reverse transcription. The
450 decrease we observed when mutating K₂₄₀ alone, to around 45 % of wt activity, is consistent with this
451 view. The characterisation by NMR of the RT-binding surface in the IN CTD, obtained using the
452 CTD₂₂₀₋₂₇₀, shows that it is made up of 9 residues (amino acids 231-258 among which K₂₄₄) that
453 strongly interact both with the RT alone (9) and with the RT/DNA complex (74). When the interaction
454 involves the complex, this surface includes 5 additional amino acids (74). Among these additional
455 residues are N₂₂₂ and K₂₄₀, which are located at one edge of the surface. It is therefore possible that
456 the nature of the residues at positions 222 and 240 affects the interaction between the CTD and
457 RT/DNA complex.

458 Concerning K₂₇₃, contradictory results have been obtained for reverse transcription of viruses
459 harbouring integrases with sequential C-ter deletions (IN₁₋₂₇₀ and IN₁₋₂₇₂) (75, 76). Furthermore, for
460 reverse transcription to occur, the genomic RNA must be encapsidated in the core of the viral particle.
461 In this sense, it has been recently shown that mutating K₂₇₃ together with R₂₆₉ (R269A/K273A mutant)
462 impairs encapsidation of the genomic RNA (6). As expected, reverse transcription in the double
463 mutant was almost abolished. Here, mutating K₂₇₃ to Q led only to a reduction of reverse transcription
464 of 30% (Figure 3B), suggesting that mutating K₂₇₃ alone is not sufficient to affect genomic
465 incorporation into the viral capsid, at least in the majority of the particles. Supporting this, an earlier
466 study (77) showed that the K273A single mutation did not affect viral replication in Jurkat cells,
467 indicating that is the specific combination of R269A/K273A mutations to be responsible for the
468 impairment of the genome encapsidation.

469 Finally, acetylation of K₂₇₃, has been previously proposed to be important for different steps of the
470 infectious cycle (78, 79). In those studies, though, the role of acetylation of K₂₇₃ was assessed by
471 simultaneously replacing K₂₆₄, K₂₆₆, and K₂₇₃, thereby not allowing to conclude on the specific
472 contribution of K₂₇₃. Here, hampering acetylation but preserving the positive charge by the K273R
473 substitution did not affect integration, indicating that the possible acetylation of K₂₇₃ had no effect on
474 integration. This observation is in line with the observation by Topper and co-workers that
475 posttranslational acetylation of the integrase CTD is dispensable for viral replication (80). Altogether,
476 the data available in the literature regarding K₂₄₀ and K₂₇₃ indicate that the effects we observed in this
477 study cannot be due to any of the already known properties of the residues of the motif.

478 The NKNK motif is strictly conserved in natural sequences of HIV-1 group M. However, we show that
479 various variants of the NKNK motif display levels of integration efficiency equivalent to the wt enzyme
480 and could therefore in principle be found in the epidemics. Their absence is indicative of purifying
481 selection occurring *in vivo*, likely exerted at a step different from integration. One possibility is the
482 implication of IN in reverse transcription, which is, in all variants, less efficient than with the NKNK
483 sequence. The existence of several alternative sequence arrangements, in the motif, with comparable
484 efficiencies of integration might therefore have constituted an asset for optimizing the acquisition of
485 additional functions, such as promoting reverse transcription.

486

487 **Materials and Methods**

488 **Plasmids and molecular cloning**

489 p8.91-MB was constructed by engineering one *MluI* and one *BspEI* restriction sites respectively 18 nt
490 downstream the 5' and 21 nt upstream the 3' of the RT-coding sequence of the pCMV Δ R8.91 (81).
491 The insertion of the two restriction sites led to three amino acids changes in the RT (E6T, T7R and
492 A554S). These modifications only slightly affected the efficiency of generation of puromycin-resistant
493 clones (see below) upon transduction with the resulting viral vector (v8.91-MB) since, in three
494 independent experiments, the number of clones obtained with p8.91-MB was 80% \pm 6% of that
495 obtained using p8.91. The p8.91-MB was employed throughout the study as positive control and was
496 used to insert the various variants of RT and IN tested. Together with the *SaII* site, present in the
497 p8.91 48 bp downstream the stop codon of *pol* gene CDS, the *MluI* and *BspEI* sites define two
498 exchangeable cassettes: one encompassing the RT coding sequence (*MluI*-RT-*BspEI*, 1680 bp) and
499 one encompassing the IN-coding sequence (*BspEI*-IN-*SaII*, 940 bp). These cassettes were used to
500 insert the various sequences of RT and IN used in the study. The plasmid used to produce the
501 genomic RNA of the viral vectors was a modified version of pSDY, previously described (82), hereafter
502 called pSRP (for pSDY-nRFP-Puro). This variant was obtained by introducing two modifications to the
503 original pSDY-dCK-Puro plasmid (82). The first one was the replacement of the sequence encoding
504 the human deoxycytidine kinase by a cassette containing the RFP fused with the N-ter 124 amino
505 acids of human histone H2B, which directs the RFP to the nucleus. The RFP was used to monitor the
506 efficiency of transfection by fluorescence microscopy. The second modification was the replacement
507 of the HIV-1 U3 sequence in the 5' LTR by that of the U3 of the Rous sarcoma virus. For the
508 generation of qPCR standard curves, two plasmids were constructed: one, called pJet-1LTR, for the
509 detection of early and late reverse transcription products, was obtained by inserting the sequence
510 encompassing the LTR and the Psi region from pSDY (82) in the pJET plasmid with the CloneJET
511 PCR Cloning Kit (Thermo Scientific, MA, USA); the second, pGenuine2LTR, has been obtained by
512 inserting a fragment of 290 bp corresponding to the unprocessed junction of U5/U3 (CAGT/ACTG
513 being the sequence of the junction 5' to 3') into the pEX-A2 plasmid (Eurofins Genomics,

514 Luxembourg). For the study with replication-competent viruses we used the pNL4.3 plamid (83) that
515 was obtained from the NIH AIDS Research and Reference Reagent Program, #114 (GeneBank
516 accession #AF324493). We replaced in this plasmid the coding sequence of NL4.3 IN CTD with those
517 of wt and mutants INA CTD, as described in Results. Chimerical integrases between primary isolates
518 from HIV-1 group M subtype A2 and HIV-1 group O RBF206, as well as mutant integrases, were
519 constructed through overlap extension PCR as previously described for the envelope gene (84).

520

521 **Cells**

522 HEK-293T cells were obtained from the American Type Culture Collection (ATCC). P4-CCR5 reporter
523 cells are HeLa CD4+ CXCR4+ CCR5+ carrying the LacZ gene under the control of the HIV-1 LTR
524 promoter (85). TZM-bl cells are a HeLa cell clone genetically engineered to express CD4, CXCR4,
525 and CCR5 and containing the Tat-responsive reporter gene for the firefly luciferase under the control
526 of the HIV-1 long terminal repeat (86). HEK-293T, P4-CCR5 and TZM-bl cells were grown in
527 Dulbecco's Modified Eagle's Medium (DMEM, Thermo Fisher, MA, USA) supplemented with 10%
528 foetal calf serum and 100 U/ml penicillin-100 mg/ml streptomycin (Thermo Fisher, MA, USA) at 37°C
529 in 5 % CO₂. CEM-SS cells are human T4-lymphoblastoid cells (87-89) and were grown in Roswell
530 Park Memorial Institute medium (RPMI) supplemented with 10% foetal calf serum and 100 U/ml
531 penicillin-100 mg/ml streptomycin (Thermo Fisher, MA, USA) at 37°C in 5 % CO₂.

532

533 **Viral strains**

534 The following primary isolates were used for this study: from HIV-1 group M, one from subtype A2
535 (GenBank accession #AF286237, named hereafter "isolate A"), one from subtype C (GenBank
536 accession #AF286224, hereafter named "isolate C"), one from CRF02_AG (GenBank accession
537 #MH351678), one from subtype B (isolate AiHo GenBank accession #MH351679, hereafter named
538 isolate B); from HIV-1 group O the primary isolate RBF 206 (GenBank accession #KU168298,
539 hereafter named "isolate O"). Isolates #AF286237, #AF286224 and #MH351678 were obtained from

540 the NIH AIDS Research and Reference Reagent Program; isolates #MH351678, #MH351679 and
541 #KU168298 were kindly provided by J.C. Plantier (CHU Rouen, France).

542

543 **Sequence alignments**

544 We used 3366 HIV-1 sequences for alignment. HIV-1 group M sequences were downloaded from the
545 Los Alamos National Laboratory (LANL) HIV sequence database and correspond to the different HIV-
546 1 group M pure subtypes: A (249 sequences), B (2450 sequences), C (450 sequences), D (121
547 sequences), G (80 sequences), H (8 sequences), J (6 sequences), K (2 sequences). We also aligned
548 49 HIV-1 group O sequences, using 26 sequences from the LANL database and the 23 sequences
549 obtained through collaboration with the Virology Unit associated to the French National HIV Reference
550 Center (Pr. J.C. Plantier). Sequence alignments were performed with CLC sequence viewer 8. The
551 sequence logo of positions 222, 240, 254 and 273 in HIV-1 group M IN was obtained with an
552 alignment of 3366 sequences of the IN CTD using WebLogo (<http://weblogo.berkeley.edu/logo.cgi>).

553

554

555 **Generation of pseudotyped viral vectors**

556 Pseudotyped lentiviral vectors were produced by co-transfection of HEK 293T cells with pHCMV-G
557 (90) encoding the VSV-G envelope protein, pSRP and p8.91-MB based plasmids with the
558 polyethylenimine method following the manufacturer's instructions (PEI, MW 25000, linear;
559 Polysciences, Warrington, PA, USA). HEK 293T were seeded at 5×10^6 per 100-mm diameter dish
560 and transfected 16-20h later. The medium was replaced 6h after transfection, and the vectors were
561 recovered from the supernatant 72h later, filtered on 0.45 μ m filters and the amount of p24 (CA) was
562 quantified by ELISA (Fujirebo Europe, Belgium).

563

564 **Western blot**

565 Western blot analysis was carried out on virions to assess the proteolytic processing of the Pr55Gag
566 polyprotein. 1.5 mL of viral supernatant was centrifuged through 20 % sucrose, and the virion pellet
567 was lysed in Laemmli buffer 1.5X. Viral proteins were separated on a Criterion™ TGX Strain-Free 4-
568 15 % gradient gel (Biorad, CA, USA) (TGS: Tris Base 0,025 M/Glycine 0,192 M/SDS 0,1 %, 150V, 45
569 min), blotted on a PVDF membrane (TGS/Ethanol 10 %, 200 mA, 1.5h) and probed with a mouse
570 monoclonal anti-CA antibody (NIH AIDS Reagent Program, #3537) to detect the viral capsid, the
571 Pr55Gag unprocessed polyprotein and CA-containing proteolytic intermediates. An anti-mouse HRP-
572 conjugated secondary antibody was used to probe the membrane previously incubated with anti-CA.
573 Membranes were incubated with ECL reagent (Thermo Fisher, MA, USA) and WB were imaged on a
574 Biorad Chemidoc Touch and analysed with the Biorad Image Lab software.

575

576 **Evaluation of reverse transcription by qPCR**

577 The viral vectors were treated with 200 U/ml of Benzonase nuclease (Sigma-Aldrich, MO, USA) in the
578 presence of 1 mM MgCl₂ for 1h at 37°C to remove non-internalized DNA. The vectors (200 ng of p24)
579 were then used to transduce 0.5 x 10⁶ HEK 293T cells by spinoculation for 2h at 32°C, 800 rcf, with
580 8 µg/mL polybrene (Sigma-Aldrich, MO, USA). After 2h, the supernatant was removed, cells were
581 resuspended in 2 mL of DMEM and plated in 6-well plates. After 30h, cells were trypsinised and
582 pelleted. Total DNA was extracted with UltraClean® Tissue & Cells DNA Isolation Kit (Ozyme,
583 France). A duplex qPCR assay (see Table S2 for primers) was used to quantify early and late reverse
584 transcription products by detecting the R-U5 and U5-Psi junctions, respectively, and another qPCR
585 (Table S2) to normalise for the quantity of cells employed in the assay (detection of β-actin exon 6
586 genomic DNA; International DNA Technologies -IDT- Belgium). All primers and probes were
587 synthesised by IDT. The qPCR assays were designed with the Taqman® hydrolysis probe technology
588 using the IDT Primers and Probes design software (IDT), with dual quencher probes (one internal
589 ZEN™ quencher and one 3' Iowa Black™ FQ quencher) (Table S2). qPCRs were performed with the
590 iTaq Universal Probes Supermix (Biorad, CA, USA) on a CFX96 (Biorad, CA, USA) thermal cycler with
591 the following cycling conditions: initial Taq activation 3 min, 95°C followed by [denaturation 10
592 sec/95°C; elongation 20 sec/55°C] x 40 cycles. Standard curves and analysis were carried out with

593 the CFX Manager (Biorad, CA, USA). DNA copy number was determined using a standard curve
594 prepared with serial dilutions of the reference plasmids pJet-1LTR and of a known number of HEK
595 293T cells.

596

597 **Evaluation of integration by puromycin assay**

598 Half a million of HEK 293T cells were transduced with a volume of viral vectors corresponding to 0.2
599 ng of p24, by spinoculation 2h at 32°C, 800 rcf, with 8µg/mL polybrene (Sigma-Aldrich, MO, USA).
600 After 2h the supernatant was removed, cells were resuspended in 7 mL of DMEM and plated in 100
601 mm diameter plates. After 30h, puromycin was added at a final concentration of 0.6 µg/mL, clones
602 were allowed to grow for 10 to 12 days and then counted. However, the number of clones depends on
603 two parameters: the efficiency of integration and the amount of pre-proviral DNAs available for
604 integration (which depends on the efficiency of reverse transcription). Therefore, we normalized the
605 number of clones observed by the amount of viral DNA generated by reverse transcription (estimated
606 by qPCR) to extrapolate the efficiency of integration. The percentage of integration efficiency for
607 sample X with respect to the control C is, thus, given by $(p_x/r_x)/(p_c/r_c) \times 100$, where r_x and r_c are the
608 amounts of late reverse transcription products, estimated by qPCR, in sample X and in control C,
609 respectively, and p_x and p_c the number of puromycin-resistant clones in sample X and in control C,
610 respectively.

611

612 **Evaluation of integration by Alu qPCR**

613 Equal amounts of total DNA extracted from transduced cells (as deduced by qPCR of β -actin exon 6
614 genomic DNA, see above) were used for the Alu PCR assay, as previously described (91). Two
615 subsequent amplification were performed. The first one, 95°C for 3 min, [95°C for 30 sec, 55°C for 30
616 sec, 72°C for 3 min30s] x15, 72°C for 7 min, using the Alu-forward primer and the Psi reverse primer,
617 allowed to amplify Alu-LTR fragments (Table S3). Samples were then diluted to 1:10 and 2 µL were
618 used for the second amplification to detect the viral LTR, as described above for the detection of the
619 R-U5 junction. The percentage of integration efficiency for sample X with respect to the control C is

620 given by $(a_x/r_x)/(a_c/r_c) \times 100$, where r_x and r_c are the amounts of late reverse transcription products,
621 estimated by qPCR, in sample X and in control C, respectively, and a_x and a_c the amounts of DNA
622 estimated by the second amplification of the Alu qPCR assay in sample X and in control C,
623 respectively.

624

625 **Quantification of two LTR circles and of circles with perfect junctions**

626 Non-internalised DNA was removed by treatment with Benzonase nuclease as for the qPCR assay
627 and 0.5×10^6 HEK 293T cells were transduced with a volume of viral vectors corresponding to 1 μ g of
628 p24 by spinoculation, as described above. After 30h, cells were trypsinised and pelleted. Total DNA
629 was extracted with UltraClean® Tissue & Cells DNA Isolation Kit. Late reverse transcription products
630 (detection of the U5-Psi junction) were quantified as described above and two qPCR assays were
631 used to quantify the 2LTRc and the quantity of 2LTR circles with a perfect palindromic junction, with a
632 primer overlapping the 2LTRc junction, as previously described (92). The qPCR assays were
633 designed and the primers and probes (Table S4) synthesised as described above. qPCRs were
634 performed as described above. Standard curves and analysis were carried out with the CFX Manager
635 (Biorad, CA, USA). Copy numbers of the different forms of viral DNA were determined with respect to
636 a standard curve prepared by serial dilutions of the pGenuine2LTR plasmid. The amount of 2LTRc for
637 sample X is normalised by the total amount of the late reverse transcription products (detection of the
638 U5-Psi junction), then it is expressed as a percentage of the amount detected for the control INA
639 D116A (indicated with D), thus giving $(2LTRc_x/r_x)/(2LTRc_D/r_D) \times 100$, where r_x and r_D are the amounts
640 of late reverse transcription products, in sample X and in control D, respectively, and $2LTRc_x$ and
641 $2LTRc_D$ the amount of 2LTR circles in sample X and in control D, respectively.

642

643 **Calculation of the efficiency of nuclear import and of 3' processing**

644 The efficiency of nuclear import was estimated as follows. The level of 2LTRc found with wt IN A was
645 0.2 with respect to that found with D116A (data from Figure 6A). The diminution observed with wt IN A
646 with respect to D116A (which was considered to produce the maximum amount of 2LTRc and was

647 therefore set at 1) was therefore 0.8 (given by $1 - 0.2$). Since the diminution of 2LTRc is proportional to
648 the efficiency of integration, for example a mutant integrating with an efficiency 0.3 that of wt IN A is
649 expected to reduce the amount of 2LTRc by $0.8 \times 0.3 = 0.24$. The amount of 2LTRc expected for that
650 mutant would therefore be given by $1 - 0.24 = 0.76$. Similarly, a mutant integrating with a higher
651 efficiency (for example 0.9 that of wt IN A) is expected to give $1 - (0.8 \times 0.9) = 0.28$ 2LTRc with respect
652 to the mutant D116A. Therefore, the formula applied to estimate the expected levels of 2LTRc with
653 respect to D116A for sample n is given by $1 - (0.8 \times a_n)$ where a_n is the level of integration measured
654 for sample n, relative to wt IN A (data from Figure 6A). The values of 2LTRc measured experimentally
655 (Table 1, line 3) are then divided by the expected ones to obtain an estimate of the relative efficiency
656 of nuclear import (Table 1, line 4).

657 The efficiency of 3' processing was calculated as follows. The ratio of perfect junctions out of the total
658 amount of 2LTRc (PJ/2LTRc) found for D116A was considered to be the maximal one and was
659 therefore assigned a value of 1 (Table 1, line 5). The proportion of PJ/2LTRc found for wt IN A (Table
660 1) was 0.54 that of D116A (Table 1, line 5). The proportion by which the pool of PJ/2LTRc found with a
661 catalytically inactive IN can be decreased by 3' processing carried out by a fully catalytic active IN is
662 therefore $1 - 0.54 = 0.46$ (Table 1, line 6). For mutant NQNK, for example, the ratio PJ/2LTRc observed
663 was 0.87 of D116A, which corresponds to a relative decrease of the PJ/2LTRc pool by 0.13 (Table 1,
664 line 6). This decrease is $0.13 / 0.46 = 0.28$ that observed for wt IN A (Table 1, line 7), providing an
665 estimate of the relative efficiency of 3' processing by this mutant with respect to wt IN A. The general
666 formula we applied to estimate the efficiency of 3' processing was therefore $(1 - r_x) / 0.46$, where 1 is the
667 proportion of PJ/2LTRc found for D116A, r_x is the ratio PJ/2LTRc observed for sample X and 0.46 is
668 the decrease in PJ/2LTRc observed for wt IN A with respect to D116A. The resulting values are
669 reported in Table 1, line 7.

670

671 **Assessment of the infectivity of replication-competent viruses**

672 As described above, the coding sequence of NL4.3 IN CTD was replaced with those of wt and
673 mutants INA CTD. Replication-competent viruses were produced as described above and equal
674 amounts of viruses were used to infect cells (TZM-bL or CEM-SS) for each sample. For estimating

675 viral replication in TZM-bL cells, 25 μ L of virus dilution were added to 10^4 cells, plated in 96 wells
676 plates in 75 μ L of culture medium. After 48h, virus replication was detected by measuring Luc reporter
677 gene expression by removing 50 μ L of culture medium from each well and adding 50 μ L of Bright Glo
678 reagent to the cells. After 2 min of incubation at room temperature to allow cell lysis, 100 μ L of cell
679 lysate were transferred to 96-well black solid plates for measurements of luminescence (RLU) using a
680 luminometer (93). For the detection of virus replication in CEM-SS cells, 0.5×10^6 CEM-SS cells/5ml
681 were infected with 1/25 virus dilution. After 5 days of culture, the percentage of infected cells were
682 detected by intracellular p24 immuno-staining and flow cytometry analysis as previously described
683 (94).

684

685 **Cloning, production, purification and crystallization**

686 The C-terminal domains (IN CTD, 220-270) of integrases NL4.3, A and A K240Q/N254K studied here
687 were cloned in the pET15b plasmid and the proteins were expressed in BL21DE3 *E. coli* cells. After
688 transformation with the IN C-ter expressing pET15b, bacteria were inoculated at an OD₆₀₀ of 0.1 in
689 one litre of LB medium supplemented with 10% (w/v) sucrose. Cultures were incubated at 37°C with
690 shaking at 220 rpm. At OD_{600nm} of 0.5, the temperature was lowered to 25°C, and shaking reduced to
691 190rpm, till the cells reached an OD_{600nm} of 0.8. IPTG was then added to a final concentration of 0.5
692 mM to induce the expression of the C-terminal domains. Cells were incubated overnight at 25°.
693 Bacteria were then collected by centrifugation.

694 For protein purification, cells were resuspended in lysis buffer (25 mM HEPES pH 8, 1 M NaCl, 10 mM
695 imidazole) in a ratio of 10 mL of buffer/gram of biomass. Roche Complete Inhibitor Cocktail tablets
696 were added at the beginning of lysis to avoid protease degradation. Cells were lysed by sonification,
697 for 1min/g of cells with pulse every 2 seconds at 40% amplitude at 4°C. The bacterial debris were
698 pelleted by ultracentrifugation at 100 000xg for 1hr at 4°C. The supernatant was then loaded on a 1
699 mL HisTrap FF Crude column (GE Healthcare) with flow rate of 1 mL/min using the AKTA FPLC.
700 Protein was eluted using a gradient up to 500 mM Imidazole in 10 column volumes. Protein
701 concentration was estimated using the Nanodrop. Subsequently, the protein sample was concentrated
702 using the Amicon Ultra 15 mL with a 3 kDa MWCO for the next purification step. A second step of

703 purification was carried out using the S75-16/60 column (GE Healthcare) in 25 mM HEPES pH 8, 1 M
704 NaCl. Samples were dialyzed into 25 mM HEPES pH 8, 150 mM NaCl for crystallization.

705 All initial crystallization conditions were determined by vapor diffusion using the TPP Labtech Mosquito
706 Crystal. 200 nL of protein (7-4 mg/mL) was mixed with 200 nL of reservoir in 2 or 3 well of a 96 well
707 MRC crystallization plate which was stored in the Formulatrix RockImager at 20°C. Screen included
708 PEGS (Hampton Research), MPD, CLASSICS, NUCLEIX (Qiagen), JCSG, WIZARDS, ANION and
709 CATION (Molecular Dimensions). Once the initial conditions were obtained, manual drops were set up
710 in Hampton Research 24 well VDX plate to optimize crystallization conditions, and to improve crystal
711 size and quality by mixing 1 μ L protein + 1 μ L reservoir and equilibrating against 500 μ L of reservoir at
712 20°C. The IN CTD NL4.3 (group M, subtype B) were obtained in a reservoir containing 0.1 M Tris pH
713 7, 0.8 M potassium sodium tartrate and 0.2 M lithium sulfate. For IN CTD A (subtype A2) and A
714 K240Q/N254K the reservoir was composed of 0.1 M MES pH 6.5 and 1M sodium malonate.

715

716 **Data collection, structure solving and refinement**

717 Data collection was performed at the Swiss Light Source (SLS, Villigen, Switzerland) on a Dectris
718 Pilatus 2M detector. After fishing, crystals were rapidly passed through a drop of fluorinated oil
719 (Fomblin® Y LVAC 14/6, average MW 2,500 from Sigma Aldrich) to prevent ice formation and directly
720 frozen on the beamline in the nitrogen stream at 100 K. X-ray diffraction images were indexed and
721 scaled with XDS (95, 96). The structures were solved by molecular replacement using PHASER (97)
722 in the PHENIX (98) program suite using the NMR HIV-1 C-ter structure (1QMC) (99) as a search
723 model for the IN CTD NL4.3 structure, which was used subsequently as a search model to solve the A
724 and A K240Q/N254K structures. The structure was then built using the AUTOBUILD program (100,
725 101) followed by several cycles of refinement using PHENIX.REFINE (102) and manual rebuilding
726 with COOT (103). Structure based sequence alignment was performed using PROMALS3D (104).
727 Structures superposition and Root Mean Square Deviations (RMSD) calculations have been
728 performed using secondary structure matching (SSM), superpose program (105) embedded in COOT
729 (103) and in the CCP4 program suite (106). The sequence alignment representation has been
730 generated by ESPript (107). Surface potential was calculated using the DELPHI web server (108) and

731 visualized with CHIMERA (109). Data collection and refinement statistics are summarized in Table S5.
732 Crystallographic structures were deposited in PDB under the identification numbers 6T6E (HIV-1 Cter,
733 PNL4.3), 6T6I (HIV-1 Cter, subtype A2) and 6T6J (HIV-1 Cter, subtype A2, mutant N254K-K240Q).

734

735 **Analysis of the surface electrostatic potential of retroviral CTDs**

736 The structures and the sequences of the C-terminal domains have been extracted from: HIV-1 A2,
737 PDB 6T6I (this publication); HIV-1 PNL4.3, PDB 6T6E (this publication); SIV, PDB 1C6V (48); MVV,
738 PDB 5LLJ (49); RSV, PDB 1C0M (50); MMTV, PDB 5D7U (51); MMLV, PDB 2M9U (52); PFV, PDB
739 4E7I (53). The structure based sequence alignment has been performed using PROMALS3D (104).
740 Structures superposition and rRoot Mean Square Deviations (RMSD) calculations have been
741 performed using secondary structure matching (SSM), superpose program (105) embedded in COOT
742 (103). The sequence alignment representation has been generated by ESPript (107). The surface
743 electrostatic potential was calculated using the DELPHI web server (108) and visualized with
744 CHIMERA (109).

745

746 **Statistical tests**

747 All statistical analyses were performed on at least three independent experiments (transfection and
748 transduction) using Prism 6 (GraphPad). For all functional tests, the values obtained for the chimeras
749 were normalized using the values obtained for parental integrase A. Student tests were used to
750 evaluate whether the normalized mean values obtained with the chimeric and mutant integrases were
751 significantly different from that obtained with the parental strain, and/or between them. For confocal
752 microscopy, unpaired t test was used for statistical analyses.

753

754

755 ***Acknowledgements***

756 The authors are grateful to Pr. J.C. Plantier for providing HIV-1 strains of subtype B, CRF02 and group
757 O, to C. Elefante for the construction of the the p8.91-MB plasmid, to J. Batisse for providing control
758 reagents, and to M. Lavigne, B. Maillot and S. Marzi for helpful discussions. The authors wish to thank
759 R. Drillien (IGBMC) for suggestions about the manuscript. The authors thank V. Olieric and the staff of
760 the Swiss Light Source synchrotron for help with data collection. The authors acknowledge the support
761 and the use of resources of the French Infrastructure for Integrated Structural Biology FRISBI ANR-
762 10-INBS-05 and of Instruct-ERIC.

763

764

765 **Table and figure legends**

766 **Table 1. Estimate of the contribution of the defects of nuclear import and 3' processing to the**
767 **efficiency of integration observed with the mutants of the NKNK motif.**

768 **Figure 1. Outline of the experimental system.** *Panel A.* Workflow used to evaluate Pr55Gag
769 processing, reverse transcription, and integration in our experimental system. VSV-pseudotyped HIV-1
770 derived vectors, produced by triple transfection, were used to transduce HEK 293T cells. Upon
771 integration, the proviral DNA will allow growth of the cellular clones in the presence of puromycin. For
772 multiplicities of infection lower than 1, the number of clones obtained is directly proportional to the
773 number of integration events. *Panel B.* schematic representation of the viral genomic RNA contained
774 in the viral vectors, transcribed from pSRP (panel A, also see Materials and Methods). R, U5 and U3,
775 viral sequences constituting the LTR; "cis-acting", viral sequences required for RNA packaging and
776 reverse transcription; EF1- α and hPGK, internal human promoters driving the expression of the
777 nuclear RFP (nRFP) and of the puromycin N-acetyl-transferase that confers resistance to puromycin
778 (Puro^R), respectively. *Panel C.* Evaluation of reverse transcription and integration in control samples,
779 compared to v8.91-MB reference vector. The results give the average values of three independent
780 experiments.

781 **Figure 2: Functionality of chimerical integrases.** *Panel A.* Alignment of CTD sequences from
782 isolates A and O, used in this study. The numbers in italic on the left and on the right of the alignment
783 indicate the beginning and the end (in amino acid) of the CTD, respectively. Only amino acids
784 divergent between the two sequences are indicated by letters. The arrows and numbers above the
785 alignment indicate the last position that, in the chimeras, was concordant with the sequence of isolate
786 A. *Panel B.* Schematic representation of the integrases studied. Integrase from isolate O is drawn at
787 the top of the panel in dark grey; integrase from isolate A is drawn at the bottom of the panel in light
788 grey. The genetic origin of the portions of the chimeras is indicated by the colour code that refers to
789 the reference isolates A and O. *Panel C.* Representative western blot obtained with an anti-CA mouse
790 monoclonal antibody. *Panel D.* Efficiency of processing of the Pr55Gag precursor, estimated by the
791 amount of CA compared to the amount of Pr55Gag precursors detected by western blot (as in panel
792 B). The results are expressed as function of the reference wt IN A, set at 100%. *Panel E.* Efficiency of

793 reverse transcription (detection of the junction U5-Psi by qPCR) expressed as function of the
794 reference wt IN A. *Panel F.* Efficiency of integration calculated with the puromycin assay, normalized
795 by the amount of total viral DNA (estimated by qPCR), expressed as function of the reference wt IN A.
796 Error bars indicate standard deviations. The results given in panels C-E are the average of 3
797 independent experiments. ** p <0.01; *** p <0.001, p values for comparison to wt IN A.

798 **Figure 3. Functionality of IN A with mutated CTD.** *Panels A-C.* Efficiency of processing of the
799 Pr55Gag precursor (panel A), of reverse transcription (panel B), and normalized efficiency of
800 integration (panel C). *Panel D.* Efficiency of processing of the Pr55Gag precursor, of reverse
801 transcription, and of normalized efficiency of integration for the K240Q/K273Q mutant. *Panel E.*
802 Efficiency of integration of the K240R/K273R mutant (NRNR in the Figure) and of wt IN A (*NKNK, set
803 at 100%). Error bars indicate standard deviations. In all panels the results are the average of 3
804 independent experiments. * p <0.05; ** p <0.01; *** p <0.001, p values for comparison to wt IN A.

805 **Figure 4. Definition of the NKNK motif and of its importance in the most widespread**
806 **phylogenetic groups of HIV-1.** *Panel A.* Efficiency of integration of various mutants of the N residues
807 of the NKNK motif and of the wt enzyme (*NKNK, set at 100%). *Panel B.* Efficiency of integration of
808 the mutant carrying the sequences of isolate O at positions 222, 240, 254 and 273 (K₂₂₂Q₂₄₀K₂₅₄Q₂₇₃,
809 KQKQ in the Figure) and of the wt enzyme (*NKNK, set at 100%). *Panel C.* Conservation logo of the
810 sequence at positions 222, 240, 254 and 273 in HIV-1 group M integrases. *Panel D.* Efficiency of
811 integration of the double mutant N240Q/K273Q (NQNQ in the Figure) of an isolate of subtype B, one
812 of subtype C and from CRF02, compared to the corresponding wt integrases, set as reference at 100
813 %. In all vectors the RT sequence had the same phylogenetic origin as IN and was replaced using the
814 *MluI-BspEI* cassette in p8.91MB, as described in Materials and Methods. Error bars indicate standard
815 deviations (standard deviations of the wt integrases of each subtype were calculated with respect to
816 reference wt IN A, used as control). The results are the average of 3 independent experiments.

817 **Figure 5. Importance of the number and position of the K residues in the N₂₂₂K₂₄₀N₂₅₄K₂₇₃ motif**
818 **of the CTD.** *Panel A.* Efficiency of integration, normalized by the amount of viral DNA, for the IN
819 mutants grouped by the number of K present at positions 222, 240, 254, 273. The composition in
820 amino acids in the four positions of the motif is given for the isolate with 0 and for the one with 4 K. For

821 clarity, only the four letters of the amino acids of the motif are represented for each mutant, omitting
822 the positions; the first letter indicates the residue at position 222, the second, position 240, the third,
823 position 254 and the fourth, position 273. *Panels B-D*. Efficiency of integration of the individual
824 mutants containing 1 (panel B), 2 (panel C), or 3 (panel D) K in the motif. In panel C, the motif
825 corresponding to the sequence of wt IN A (reference set at 100%) is indicated by an asterisk. Error
826 bars indicate standard deviations. The results are the average of 4 independent experiments. * p
827 <0.05 ; ** $p <0.01$; *** $p <0.001$, p values for comparison to wt IN A. *Panels E-F*, importance of the
828 NKNK motif in replication-competent viruses. *Panel E*. Scheme of the portion coding for the integrase
829 in the various viruses. Drawn in grey are the parts derived from the NL4.3 sequences, in white those
830 from isolate A. The black bars indicate positions 222, 240, 254 and 273 from left to right; the amino
831 acid found for each mutant in each of these four positions is indicated above the bars. *Panel F*.
832 Infectivity of the viruses shown in panel E (except wt NL4.3 that is used as reference, set at 100%).
833 The results are given in grey for CEM-SS and in black for TZM-bL cells. Error bars indicate standard
834 deviations with respect to the reference wt pNL4-3. The results are the average of 2 independent
835 experiments.

836 **Figure 6. Amount of 2LTRc (panel A) and of ratio of PJ/2LTRc (panel B) in the mutants deprived**
837 **of one or both K of the NKNK motif.** The motif corresponding to the sequence of wt IN A (reference
838 set at 100%) is indicated by an asterisk. Error bars indicate standard deviations. Above the plot are
839 given the p values for the comparisons of the different samples with respect to wt IN A or to the
840 integration-deficient mutant IN A D116A (* $p <0.05$; ** $p <0.01$; *** $p <0.001$). The number of
841 independent experiments performed for each sample (n) is also given.

842 **Figure 7. Structural analysis of the NKNK motif.** *Panels A and B*. Side view (A) and top view (B) of
843 the ribbon representation of the crystal structure of CTD A. The positions of residues N222, K240 and
844 N254 are represented with sticks as well as the position of the I234, the only different residue between
845 the CTD of IN A and IN NL4.3. *Panels C (side view) and D (top view)* are the surface electrostatic
846 potential representation of CTD A. In red, negative potential; in blue, positive potential and in white
847 neutral regions. Circled in yellow is the region with large differences in the mutant structures (see
848 below *Panel H-M*). *Panel E*. Superposition of the CTDs of IN A and IN A NQKK (chain A, B and C).
849 The mutation N254K induces a displacement of the K236 side chain (white arrows) disturbing the

850 structure of the 235-237 region. *Panels F and G.* Side view (F) and top view (G) of the superposition of
851 the three molecules in the asymmetric unit of the NQKK CTD. The position of the residue N222,
852 K240Q mutation and N254K mutation are represented with sticks as well as the position of the I234.
853 *Panels H, J and L* (side view) and *I, K and M* (top view) are the surface electrostatic potential
854 representation of NQKK CTD chain A (H, I), chain B (J, K) and chain C (L, M). In red, negative
855 potential; in blue, positive potential and in white neutral regions. The regions with large differences are
856 circled in yellow.

857 **Figure 8. Analysis of the surface electrostatic potential in the C-ter of retroviral integrases.** The
858 structures and the sequences of the C-terminal domains have been extracted from: HIV-1 A2, PDB
859 6T6I (this publication); HIV-1 PNL4.3, PDB 6T6E (this publication); SIV, PDB 1C6V (48); MVV, PDB
860 5LLJ (49); RSV, PDB 1C0M (50); MMTV, PDB 5D7U (51); MMLV, PDB 2M9U (52); PFV, PDB 4E7I
861 (53). *Panel A.* Superposition of the structure of the integrase C-terminal domains of four lentiviruses
862 (HIV-1 A2, pink; HIV-1 pNL4.3, orange; Simian Immunodeficiency Virus, SIV, khaki; Maedi-Visna
863 virus, MVV, cyan), of an α retrovirus (the Rous Sarcoma Virus, RSV, blue), of a β retrovirus (the
864 Mouse Mammary Tumor Virus, MMTV, sky blue), of a γ retrovirus (the Moloney Murine Leukemia
865 Virus, MMLV, purple), and of a spumaretrovirus (the Prototype Foamy Virus, PFV, green). *Panel B.*
866 Structure-based sequences alignment of the integrases C-terminal domains. Sequence numbering
867 corresponds to the HIV-1 A2 integrase sequence. Secondary structures from HIV-1 A2 are
868 represented (TT: β -Turn, β 1 to β 5: β -sheets, η 1: 3_{10} -helix). Residues framed in blue: Position in the
869 alignment of the three first amino acids from the NKNK motif. Red background: 100% identity in the
870 sequence alignment. Yellow background: % of equivalent residues > 70% (considering their physical-
871 chemical properties), equivalent residues are depicted in bold. *Panels C-L.* Surface electrostatic
872 potential representation of integrases from several retroviruses. The surface corresponding to that
873 delimited by the NKNK motif in HIV-1 M (panels C, D and E) is circled in yellow. *Panel C,* ribbon
874 representation of HIV-1 A2 C-terminal domain structure. The amino acids belonging to the NKNK motif
875 are represented in sticks. *Panels D-G.* Surface potential representation of the C-terminal domain
876 structures of four lentiviral integrases: HIV-1 A2 (panel D), HIV-1 pNL4.3 (panel E), Simian
877 Immunodeficiency Virus (SIV) (panel F) and Maedi-Visna virus (MVV) (panel G). *Panel H.* Ribbon
878 representation of the structure of Rous Sarcoma Virus (RSV) C-terminal domain. The surface circled

879 in yellow corresponds to that delimited by the NKNK motif in HIV-1 M after superposition of the
880 structures. The amino acids corresponding to the motif in the structure-based alignment are shown as
881 sticks. *Panel I.* Surface potential representation of the C-terminal domain structure of an α retrovirus,
882 the Rous Sarcoma Virus (RSV). *Panel J.* Surface potential representation of the C-terminal domain
883 structure of a β retrovirus, the Mouse Mammary Tumor Virus (MMTV). *Panel K.* Surface potential
884 representation of the C-terminal domain structure of a γ retrovirus, the Moloney Murine Leukemia
885 Virus (MMLV). *Panel L.* Surface potential representation of the C-terminal domain structure of a
886 spumaretrovirus, the Prototype Foamy Virus (PFV). Negative potential is in red, neutral in white and
887 positive potential is in blue.

888
889
890
891

References

- 892 1. **Pauza CD.** 1990. Two bases are deleted from the termini of HIV-1 linear DNA during
893 integrative recombination. *Virology* **179**:886–889.
- 894 2. **Engelman A, Mizuuchi K, Craigie R.** 1991. HIV-1 DNA Integration - Mechanism of Viral-
895 Dna Cleavage and Dna Strand Transfer. *Cell* **67**:1211–1221.
- 896 3. **Engelman A, Englund G, Orenstein JM, Martin MA, Craigie R.** 1995. Multiple Effects of
897 Mutations in Human-Immunodeficiency-Virus Type-1 Integrase on Viral Replication. *J Virol*
898 **69**:2729–2736.
- 899 4. **Bukovsky A, Gottlinger H.** 1996. Lack of integrase can markedly affect human
900 immunodeficiency virus type 1 particle production in the presence of an active viral protease.
901 *J Virol* **70**:6820–6825.
- 902 5. **Hoyte AC, Jamin AV, Koneru PC, Kobe MJ, Larue RC, Fuchs JR, Engelman AN,**
903 **Kvaratskhelia M.** 2017. Resistance to pyridine-based inhibitor KF116 reveals an unexpected
904 role of integrase in HIV-1 Gag-Pol polyprotein proteolytic processing. *J Biol Chem*
905 **292**:19814–19825.
- 906 6. **Kessl JJ, Kutluay SB, Townsend D, Rebensburg S, Slaughter A, Larue RC, Shkriabai N,**
907 **Bakouche N, Fuchs JR, Bieniasz PD, Kvaratskhelia M.** 2016. HIV-1 Integrase binds the
908 viral RNA genome and is essential during virion morphogenesis. *Cell* **166**:1257–1268.e12.
- 909 7. **Zhu K, Dobard C, Chow SA.** 2004. Requirement for integrase during reverse transcription of
910 human immunodeficiency virus type 1 and the effect of cysteine mutations of integrase on its
911 interactions with reverse transcriptase. *J Virol* **78**:5045–5055.
- 912 8. **Dobard CW, Briones MS, Chow SA.** 2007. Molecular mechanisms by which human
913 immunodeficiency virus type 1 integrase stimulates the early steps of reverse transcription. *J*
914 *Virol* **81**:10037–10046.

- 915 9. **Wilkinson TA, Januszyk K, Phillips ML, Tekeste SS, Zhang M, Miller JT, Le Grice SFJ,**
916 **Clubb RT, Chow SA.** 2009. Identifying and characterizing a functional HIV-1 reverse
917 transcriptase-binding site on integrase. *J Biol Chem* **284**:7931–7939.
- 918 10. **Emiliani S, Mousnier A, Busschots K, Maroun M, Van Maele B, Tempé D,**
919 **Vandekerckhove L, Moisant F, Ben-Slama L, Witvrouw M, Christ F, Rain J-C,**
920 **Dargemont C, Debyser Z, Benarous R.** 2005. Integrase mutants defective for interaction
921 with LEDGF/p75 are impaired in chromosome tethering and HIV-1 replication. *J Biol Chem*
922 **280**:25517–25523.
- 923 11. **Bukrinsky MI, Sharova N, Dempsey MP, Stanwick TL, Bukrinskaya AG, Haggerty S,**
924 **Stevenson M.** 1992. Active nuclear import of human immunodeficiency virus type 1
925 preintegration complexes. *PNAS* **89**:6580–6584.
- 926 12. **Mattaj IW, Englmeier L.** 1998. Nucleocytoplasmic transport: the soluble phase. *Annu Rev*
927 *Biochem* **67**:265–306.
- 928 13. **Yamashita M, Emerman M.** 2004. Capsid is a dominant determinant of retrovirus infectivity
929 in nondividing cells. *J Virol* **78**:5670–5678.
- 930 14. **Yamashita M, Perez O, Hope TJ, Emerman M.** 2007. Evidence for Direct Involvement of
931 the Capsid Protein in HIV Infection of Nondividing Cells. *PLoS Pathog* **3**:e156.
- 932 15. **Di Nunzio F, Danckaert A, Fricke T, Perez P, Fernandez J, Perret E, Roux P, Shorte S,**
933 **Charneau P, Diaz-Griffero F, Arhel NJ.** 2012. Human nucleoporins promote HIV-1 docking
934 at the nuclear pore, nuclear import and integration. *PLoS ONE* **7**:e46037.
- 935 16. **Di Nunzio F, Fricke T, Miccio A, Valle-Casuso JC, Perez P, Souque P, Rizzi E,**
936 **Severgnini M, Mavilio F, Charneau P, Diaz-Griffero F.** 2013. Nup153 and Nup98 bind the
937 HIV-1 core and contribute to the early steps of HIV-1 replication. *Virology* **440**:8–18.
- 938 17. **Matreyek KA, Engelman A.** 2011. The requirement for nucleoporin NUP153 during human
939 immunodeficiency virus type 1 infection is determined by the viral capsid. *J Virol* **85**:7818–
940 7827.
- 941 18. **Krishnan L, Matreyek KA, Oztop I, Lee K, Tipper CH, Li X, Dar MJ, KewalRamani VN,**
942 **Engelman A.** 2010. The requirement for cellular transportin 3 (TNPO3 or TRN-SR2) during
943 infection maps to human immunodeficiency virus type 1 capsid and not integrase. *J Virol*
944 **84**:397–406.
- 945 19. **Cribier A, Ségéral E, Delelis O, Parissi V, Simon A, Ruff M, Benarous R, Emiliani S.**
946 2011. Mutations affecting interaction of integrase with TNPO3 do not prevent HIV-1 cDNA
947 nuclear import. *Retrovirology* **8**:104.
- 948 20. **Gallay P, Hope T, Chin D, Trono D.** 1997. HIV-1 infection of nondividing cells through the
949 recognition of integrase by the importin/karyopherin pathway. *PNAS* **94**:9825–9830.
- 950 21. **Bouyac-Bertoia M, Dvorin JD, Fouchier RAM, Jenkins Y, Meyer BE, Wu LI, Emerman M,**
951 **Malim MH.** 2001. HIV-1 infection requires a functional integrase NLS. *Molecular Cell* **7**:1025–
952 1035.
- 953 22. **Hearps AC, Jans DA.** 2006. HIV-1 integrase is capable of targeting DNA to the nucleus via
954 an Importin α/β -dependent mechanism. *Biochemical Journal* **398**:475–484.
- 955 23. **Ao Z, Huang G, Yao H, Xu Z, Labine M, Cochrane AW, Yao X.** 2007. Interaction of human
956 immunodeficiency virus type 1 integrase with cellular nuclear import receptor importin 7 and
957 its impact on viral replication. *J Biol Chem* **282**:13456–13467.

- 958 24. **Ao Z, Jayappa KD, Wang B, Zheng Y, Kung S, Rassart E, Depping R, Kohler M, Cohen**
959 **EA, Yao X.** 2010. Importin $\alpha 3$ Interacts with HIV-1 Integrase and Contributes to HIV-1
960 Nuclear Import and Replication. *J Virol* **84**:8650–8663.
- 961 25. **Woodward CL, Prakobwanakit S, Mosessian S, Chow SA.** 2009. Integrase interacts with
962 nucleoporin NUP153 to mediate the nuclear import of human immunodeficiency virus type 1.
963 *J Virol* **83**:6522–6533.
- 964 26. **Ao Z, Jayappa KD, Wang B, Zheng Y, Wang X, Peng J, Yao X.** 2012. Contribution of host
965 nucleoporin 62 in HIV-1 integrase chromatin association and viral DNA integration. *J Biol*
966 *Chem* **287**:10544–10555.
- 967 27. **Larue R, Gupta K, Wuensch C, Shkriabai N, Kessi JJ, Danhart E, Feng L, Taltynov O,**
968 **Christ F, Van Duyne GD, Debyser Z, Foster MP, Kvaratskhelia M.** 2012. Interaction of the
969 HIV-1 intasome with transportin 3 protein (TNPO3 or TRN-SR2). *J Biol Chem* **287**:34044–
970 34058.
- 971 28. **De Houwer S, Demeulemeester J, Thys W, Rocha S, Dirix L, Gijssbers R, Christ F,**
972 **Debyser Z.** 2014. The HIV-1 integrase mutant R263A/K264A is 2-fold defective for TRN-SR2
973 binding and viral nuclear import. *J Biol Chem* **289**:25351–25361.
- 974 29. **Christ F, Thys W, De Rijck J, Gijssbers R, Albanese A, Arosio D, Emiliani S, Rain J-C,**
975 **Benarous R, Cereseto A, Debyser Z.** 2008. Transportin-SR2 imports HIV into the nucleus.
976 *Curr Biol* **18**:1192–1202.
- 977 30. **Jayappa KD, Ao Z, Yang M, Wang J, Yao X.** 2011. Identification of critical motifs within HIV-
978 1 integrase required for importin $\alpha 3$ interaction and viral cDNA nuclear import. *J Mol Biol*
979 **410**:847–862.
- 980 31. **Passos DO, Li M, Yang R, Rebensburg SV, Ghirlando R, Jeon Y, Shkriabai N,**
981 **Kvaratskhelia M, Craigie R, Lyumkis D.** 2017. Cryo-EM structures and atomic model of the
982 HIV-1 strand transfer complex intasome. *Science* **355**:89–92.
- 983 32. **Michel F, Crucifix C, Granger F, Eiler S, Mouscadet J-F, Korolev S, Agapkina J,**
984 **Ziganshin R, Gottikh M, Nazabal A, Emiliani S, Benarous R, Moras D, Schultz P, Ruff M.**
985 2009. Structural basis for HIV-1 DNA integration in the human genome, role of the
986 LEDGF/P75 cofactor. *EMBO J* **28**:980–991.
- 987 33. **Craigie R, Bushman FD.** 2012. HIV DNA integration. *Cold Spring Harb Perspect Med*
988 **2**:a006890–a006890.
- 989 34. **Delelis O, Carayon K, Saïb A, Deprez E, Mouscadet J-F.** 2008. Integrase and integration:
990 biochemical activities of HIV-1 integrase. *Retrovirology* **5**:114.
- 991 35. **Zheng R, Jenkins TM, Craigie R.** 1996. Zinc folds the N-terminal domain of HIV-1 integrase,
992 promotes multimerization, and enhances catalytic activity. *Proc Natl Acad Sci USA*
993 **93**:13659–13664.
- 994 36. **Eijkelenboom AP, van den Ent FM, Vos A, Doreleijers JF, Hård K, Tullius TD, Plasterk**
995 **RH, Kaptein R, Boelens R.** 1997. The solution structure of the amino-terminal HHCC
996 domain of HIV-2 integrase: a three-helix bundle stabilized by zinc. *Curr Biol* **7**:739–746.
- 997 37. **Busschots K, Vercammen J, Emiliani S, Benarous R, Engelborghs Y, Christ F, Debyser**
998 **Z.** 2005. The interaction of LEDGF/p75 with integrase is lentivirus-specific and promotes
999 DNA binding. *J Biol Chem* **280**:17841–17847.
- 1000 38. **Heuer TS, Brown PO.** 1997. Mapping features of HIV-1 integrase near selected sites on viral
1001 and target DNA molecules in an active enzyme-DNA complex by photo-cross-linking.
1002 *Biochemistry* **36**:10655–10665.

- 1003 39. **Esposito D, Craigie R.** 1998. Sequence specificity of viral end DNA binding by HIV-1
1004 integrase reveals critical regions for protein-DNA interaction. *EMBO J* **17**:5832–5843.
- 1005 40. **Chen A, Weber IT, Harrison RW, Leis J.** 2006. Identification of amino acids in HIV-1 and
1006 avian sarcoma virus integrase subsites required for specific recognition of the long terminal
1007 repeat Ends. *J Biol Chem* **281**:4173–4182.
- 1008 41. **Engelman A, Hickman AB, Craigie R.** 1994. The core and carboxyl-terminal domains of the
1009 integrase protein of human immunodeficiency virus type 1 each contribute to nonspecific
1010 DNA binding. *J Virol* **68**:5911–5917.
- 1011 42. **Lutzke RA, Vink C, Plasterk RH.** 1994. Characterization of the minimal DNA-binding
1012 domain of the HIV integrase protein. *Nucleic Acids Res* **22**:4125–4131.
- 1013 43. **Cannon PM, Byles ED, Kingsman SM, Kingsman AJ.** 1996. Conserved sequences in the
1014 carboxyl terminus of integrase that are essential for human immunodeficiency virus type 1
1015 replication. *J Virol* **70**:651–657.
- 1016 44. **Larder BA, Purifoy DJ, Powell KL, Darby G.** 1987. Site-specific mutagenesis of AIDS virus
1017 reverse transcriptase. *Nature* **327**:716–717.
- 1018 45. **Mandal D, Feng Z, Stoltzfus CM.** 2008. Gag-processing defect of human immunodeficiency
1019 virus type 1 integrase E246 and G247 mutants is caused by activation of an overlapping 5'
1020 splice site. *J Virol* **82**:1600–1604.
- 1021 46. **Hemelaar J, Gouws E, Ghys PD, Osmanov S, WHO-UNAIDS Network for HIV Isolation
1022 and Characterisation.** 2011. Global trends in molecular epidemiology of HIV-1 during 2000-
1023 2007. *AIDS* **25**:679–689.
- 1024 47. **Sloan RD, Wainberg MA.** 2011. The role of unintegrated DNA in HIV infection. *Retrovirology*
1025 **8**:52.
- 1026 48. **Chen Z, Yan Y, Munshi S, Li Y, Zugay-Murphy J, Xu B, Witmer M, Felock P, Wolfe A,
1027 Sardana V, Emini EA, Hazuda D, Kuo LC.** 2000. X-ray structure of simian
1028 immunodeficiency virus integrase containing the core and C-terminal domain (residues 50-
1029 293): an initial glance of the viral DNA binding platform. *J Mol Biol* **296**:521–533.
- 1030 49. **Ballandras-Colas A, Maskell DP, Serrao E, Locke J, Swuec P, Jónsson SR, Kotecha A,
1031 Cook NJ, Pye VE, Taylor IA, Andrésdóttir V, Engelman AN, Costa A, Cherepanov P.**
1032 2017. A supramolecular assembly mediates lentiviral DNA integration. *Science* **355**:93–95.
- 1033 50. **Yang ZN, Mueser TC, Bushman FD, Hyde CC.** 2000. Crystal structure of an active two-
1034 domain derivative of Rous sarcoma virus integrase. *J Mol Biol* **296**:535–548.
- 1035 51. **Ballandras-Colas A, Brown M, Cook NJ, Dewdney TG, Demeler B, Cherepanov P,
1036 Lyumkis D, Engelman AN.** 2016. Cryo-EM reveals a novel octameric integrase structure for
1037 betaretroviral intasome function. *Nature* **530**:358–361.
- 1038 52. **Aiyer S, Swapna GVT, Malani N, Aramini JM, Schneider WM, Plumb MR, Ghanem M,
1039 Larue RC, Sharma A, Studamire B, Kvaratskhelia M, Bushman FD, Montelione GT,
1040 Roth MJ.** 2014. Altering murine leukemia virus integration through disruption of the integrase
1041 and BET protein family interaction. *Nucleic Acids Res* **42**:5917–5928.
- 1042 53. **Hare S, Maertens GN, Cherepanov P.** 2012. 3'-processing and strand transfer catalysed by
1043 retroviral integrase in crystallo. *EMBO J* **31**:3020–3028.

- 1044 54. **Luscombe NM, Laskowski RA, Thornton JM.** 2001. Amino acid-base interactions: a three-
1045 dimensional analysis of protein-DNA interactions at an atomic level. *Nucleic Acids Res*
1046 **29**:2860–2874.
- 1047 55. **Zhou H-X, Pang X.** 2018. Electrostatic interactions in protein structure, folding, binding, and
1048 condensation. *Chem Rev* **118**:1691–1741.
- 1049 56. **Karlin S, Zhu ZY.** 1996. Characterizations of diverse residue clusters in protein three-
1050 dimensional structures. *PNAS* **93**:8344–8349.
- 1051 57. **Karlin S, Brendel V.** 1988. Charge configurations in viral proteins. *PNAS* **85**:9396–9400.
- 1052 58. **Parker MS, Balasubramaniam A, Parker SL.** 2012. On the segregation of protein ionic
1053 residues by charge type. *Amino Acids* **43**:2231–2247.
- 1054 59. **Kharrat N, Belmabrouk S, Abdelhedi R, Benmarzoug R, Assidi M, Qahtani AI MH, Rebai**
1055 **A.** 2016. Screening for clusters of charge in human virus proteomes. *BMC Genomics*
1056 **17**:758–19.
- 1057 60. **Dostálková A, Kaufman F, Křížová I, Kultová A, Strohalmová K, Hadravová R, Ruml T,**
1058 **Rumlová M.** 2018. Mutations in the basic region of the Mason-Pfizer Monkey virus
1059 nucleocapsid protein affect reverse transcription, genomic RNA packaging, and the virus
1060 assembly site. *J Virol* **92**:5439.
- 1061 61. **Hamann MV, Müllers E, Reh J, Stanke N, Effantin G, Weissenhorn W, Lindemann D.**
1062 2014. The cooperative function of arginine residues in the Prototype Foamy Virus Gag C-
1063 terminus mediates viral and cellular RNA encapsidation. *Retrovirology*, 6 ed. **11**:87–17.
- 1064 62. **Heyrana KJ, Goh BC, Perilla JR, Nguyen T-LN, England MR, Bewley MC, Schulten K,**
1065 **Craven RC.** 2016. Contributions of charged residues in structurally dynamic capsid surface
1066 loops to Rous Sarcoma virus assembly. *J Virol* **90**:5700–5714.
- 1067 63. **Yenamandra SP, Sompallae R, Klein G, Kashuba E.** 2009. Comparative analysis of the
1068 Epstein-Barr virus encoded nuclear proteins of EBNA-3 family. *Comput Biol Med* **39**:1036–
1069 1042.
- 1070 64. **Roller RJ, Bjerke SL, Haugo AC, Hanson S.** 2010. Analysis of a charge cluster mutation of
1071 herpes simplex virus type 1 UL34 and its extragenic suppressor suggests a novel interaction
1072 between pUL34 and pUL31 that is necessary for membrane curvature around capsids. *J Virol*
1073 **84**:3921–3934.
- 1074 65. **Zhang J, Webb DM.** 2004. Rapid evolution of primate antiviral enzyme APOBEC3G. *Hum*
1075 *Mol Genet* **13**:1785–1791.
- 1076 66. **Vasudev PG, Banerjee M, Ramakrishnan C, Balaram P.** 2012. Asparagine and glutamine
1077 differ in their propensities to form specific side chain-backbone hydrogen bonded motifs in
1078 proteins. *Proteins* **80**:991–1002.
- 1079 67. **Weichenberger CX, Sippl MJ.** 2006. Self-consistent assignment of asparagine and
1080 glutamine amide rotamers in protein crystal structures. *Structure* **14**:967–972.
- 1081 68. **Métifiot M, Johnson BC, Kiselev E, Marler L, Zhao XZ, Burke TR, Marchand C, Hughes**
1082 **SH, Pommier Y.** 2016. Selectivity for strand-transfer over 3'-processing and susceptibility to
1083 clinical resistance of HIV-1 integrase inhibitors are driven by key enzyme-DNA interactions in
1084 the active site. *Nucleic Acids Res* **44**:6896–6906.

- 1085 69. **Johnson AA, Santos W, Pais GCG, Marchand C, Amin R, Burke TR, Verdine G,**
1086 **Pommier Y.** 2006. Integration requires a specific interaction of the donor DNA terminal 5'-
1087 cytosine with glutamine 148 of the HIV-1 integrase flexible loop. *J Biol Chem* **281**:461–467.
- 1088 70. **Guiot E, Carayon K, Delelis O, Simon F, Tauc P, Zubin E, Gottikh M, Mouscadet J-F,**
1089 **Brochon J-C, Deprez E.** 2006. Relationship between the oligomeric status of HIV-1
1090 integrase on DNA and enzymatic activity. *J Biol Chem* **281**:22707–22719.
- 1091 71. **Hehl EA, Joshi P, Kalpana GV, Prasad VR.** 2004. Interaction between human
1092 immunodeficiency virus type 1 reverse transcriptase and integrase proteins. *J Virol* **78**:5056–
1093 5067.
- 1094 72. **Ao Z, Fowke KR, Cohen EA, Yao X.** 2005. Contribution of the C-terminal tri-lysine regions
1095 of human immunodeficiency virus type 1 integrase for efficient reverse transcription and viral
1096 DNA nuclear import. *Retrovirology* **2**:62.
- 1097 73. **Williams KL, Zhang Y, Shkriabai N, Karki RG, Nicklaus MC, Kotrikadze N, Hess S, Le**
1098 **Grice SFJ, Craigie R, Pathak VK, Kvaratskhelia M.** 2005. Mass spectrometric analysis of
1099 the HIV-1 integrase-pyridoxal 5'-phosphate complex reveals a new binding site for a
1100 nucleotide inhibitor. *J Biol Chem* **280**:7949–7955.
- 1101 74. **Tekeste SS, Wilkinson TA, Weiner EM, Xu X, Miller JT, Le Grice SFJ, Clubb RT, Chow**
1102 **SA.** 2015. Interaction between Reverse Transcriptase and Integrase Is Required for Reverse
1103 Transcription during HIV-1 Replication. *J Virol* **89**:12058–12069.
- 1104 75. **Dar MJ, Monel B, Krishnan L, Shun M-C, Di Nunzio F, Helland DE, Engelman A.** 2009.
1105 Biochemical and virological analysis of the 18-residue C-terminal tail of HIV-1 integrase.
1106 *Retrovirology* **6**:94.
- 1107 76. **Mohammed KD, Topper MB, Muesing MA.** 2011. Sequential deletion of the integrase
1108 (Gag-Pol) carboxyl terminus reveals distinct phenotypic classes of defective HIV-1. *J Virol*
1109 **85**:4654–4666.
- 1110 77. **Lu R, Ghory HZ, Engelman A.** 2005. Genetic analyses of conserved residues in the
1111 carboxyl-terminal domain of human immunodeficiency virus type 1 integrase. *J Virol*
1112 **79**:10356–10368.
- 1113 78. **Cereseto A, Manganaro L, Gutierrez MI, Terreni M, Fittipaldi A, Lusic M, Marcello A,**
1114 **Giacca M.** 2005. Acetylation of HIV-1 integrase by p300 regulates viral integration. *EMBO J*
1115 **24**:3070–3081.
- 1116 79. **Terreni M, Valentini P, Liverani V, Gutierrez MI, Di Primio C, Di Fenza A, Tozzini V,**
1117 **Allouch A, Albanese A, Giacca M, Cereseto A.** 2010. GCN5-dependent acetylation of HIV-
1118 1 integrase enhances viral integration. *Retrovirology* **7**:18.
- 1119 80. **Topper M, Luo Y, Zhadina M, Mohammed K, Smith L, Muesing MA.** 2007.
1120 Posttranslational acetylation of the human immunodeficiency virus type 1 integrase carboxyl-
1121 terminal domain is dispensable for viral replication. *J Virol* **81**:3012–3017.
- 1122 81. **Zufferey R, Nagy D, Mandel RJ, Naldini L, Trono D.** 1997. Multiply attenuated lentiviral
1123 vector achieves efficient gene delivery in vivo. *Nat Biotechnol* **15**:871–875.
- 1124 82. **Rossolillo P, Winter F, Simon-Loriere E, Gallois-Montbrun S, Negroni M.** 2012.
1125 Retroevolution: HIV-driven evolution of cellular genes and improvement of anticancer drug
1126 activation. *PLoS Genet* **8**:e1002904.
- 1127 83. **Adachi A, Gendelman HE, Koenig S, Folks T, Willey R, Rabson A, Martin MA.** 1986.
1128 Production of acquired immunodeficiency syndrome-associated retrovirus in human and
1129 nonhuman cells transfected with an infectious molecular clone. *J Virol* **59**:284–291.

- 1130 84. **Gasser R, Hamoudi M, Pellicciotta M, Zhou Z, Visdeloup C, Colin P, Braibant M, Lagane**
1131 **B, Negroni M.** 2016. Buffering deleterious polymorphisms in highly constrained parts of HIV-
1132 1 envelope by flexible regions. *Retrovirology* **13**:50.
- 1133 85. **Charneau P, Mirambeau G, Roux P, Paulous S, Buc H, Clavel F.** 1994. HIV-1 reverse
1134 transcription. A termination step at the center of the genome. *J Mol Biol* **241**:651–662.
- 1135 86. **Rosen CA, Sodroski JG, Campbell K, Haseltine WA.** 1986. Construction of recombinant
1136 murine retroviruses that express the human T-cell leukemia virus type II and human T-cell
1137 lymphotropic virus type III trans activator genes. *J Virol* **57**:379–384.
- 1138 87. **Foley Ge, Lazarus H, Farber S, Uzman Bg, Boone Ba, Mccarthy Re.** 1965. Continuous
1139 culture of human lymphoblasts from peripheral blood of a child with acute leukemia. *Cancer*
1140 **18**:522–529.
- 1141 88. **Nara PL, Fischinger PJ.** 1988. Quantitative infectivity assay for HIV-1 and-2. *Nature*
1142 **332**:469–470.
- 1143 89. **Nara PL, Hatch WC, Dunlop NM, Robey WG, Arthur LO, Gonda MA, Fischinger PJ.**
1144 1987. Simple, rapid, quantitative, syncytium-forming microassay for the detection of human
1145 immunodeficiency virus neutralizing antibody. *AIDS Res Hum Retroviruses* **3**:283–302.
- 1146 90. **Naldini L, Blömer U, Gallay P, Ory D, Mulligan R, Gage FH, Verma IM, Trono D.** 1996. In
1147 vivo gene delivery and stable transduction of nondividing cells by a lentiviral vector. *Science*
1148 **272**:263–267.
- 1149 91. **Vozzolo L, Loh B, Gane PJ, Tribak M, Zhou L, Anderson I, Nyakatura E, Jenner RG,**
1150 **Selwood D, Fassati A.** 2010. Gyrase B inhibitor impairs HIV-1 replication by targeting Hsp90
1151 and the capsid protein. *J Biol Chem* **285**:39314–39328.
- 1152 92. **De Iaco A, Santoni F, Vannier A, Guipponi M, Antonarakis S, Luban J.** 2013. TNPO3
1153 protects HIV-1 replication from CPSF6-mediated capsid stabilization in the host cell
1154 cytoplasm. *Retrovirology* **10**:20.
- 1155 93. **Sarzotti-Kelsoe M, Bailer RT, Turk E, Lin C-L, Bilaska M, Greene KM, Gao H, Todd CA,**
1156 **Ozaki DA, Seaman MS, Mascola JR, Montefiori DC.** 2014. Optimization and validation of
1157 the TZM-bl assay for standardized assessments of neutralizing antibodies against HIV-1. *J*
1158 *Immunol Methods* **409**:131–146.
- 1159 94. **Lederle A, Su B, Holl V, Penichon J, Schmidt S, Decoville T, Laumond G, Moog C.**
1160 2014. Neutralizing antibodies inhibit HIV-1 infection of plasmacytoid dendritic cells by an
1161 FcγRIIIa independent mechanism and do not diminish cytokines production. *Sci Rep* **4**:5845.
- 1162 95. **Kabsch W.** 2010. Integration, scaling, space-group assignment and post-refinement. *Acta*
1163 *Crystallogr D Biol Crystallogr* **66**:133–144.
- 1164 96. **Kabsch W.** 2010. XDS. *Acta Crystallogr D Biol Crystallogr* **66**:125–132.
- 1165 97. **McCoy AJ, Grosse-Kunstleve RW, Adams PD, Winn MD, Storoni LC, Read RJ.** 2007.
1166 Phaser crystallographic software. *J Appl Crystallogr* **40**:658–674.
- 1167 98. **Adams PD, Afonine PV, Bunkóczi G, Chen VB, Davis IW, Echols N, Headd JJ, Hung L-**
1168 **W, Kapral GJ, Grosse-Kunstleve RW, McCoy AJ, Moriarty NW, Oeffner R, Read RJ,**
1169 **Richardson DC, Richardson JS, Terwilliger TC, Zwart PH.** 2010. PHENIX: a
1170 comprehensive Python-based system for macromolecular structure solution. *Acta Crystallogr*
1171 *D Biol Crystallogr* **66**:213–221.

- 1172 99. **Eijkelenboom AP, Sprangers R, Hård K, Puras Lutzke RA, Plasterk RH, Boelens R,**
1173 **Kaptein R.** 1999. Refined solution structure of the C-terminal DNA-binding domain of human
1174 immunovirus-1 integrase. *Proteins* **36**:556–564.
- 1175 100. **Terwilliger TC, Grosse-Kunstleve RW, Afonine PV, Moriarty NW, Adams PD, Read RJ,**
1176 **Zwart PH, Hung L-W.** 2008. Iterative-build OMIT maps: map improvement by iterative model
1177 building and refinement without model bias. *Acta Crystallogr D Biol Crystallogr* **64**:515–524.
- 1178 101. **Terwilliger TC, Grosse-Kunstleve RW, Afonine PV, Moriarty NW, Zwart PH, Hung L-W,**
1179 **Read RJ, Adams PD.** 2008. Iterative model building, structure refinement and density
1180 modification with the PHENIX AutoBuild wizard. *Acta Crystallogr D Biol Crystallogr* **64**:61–69.
- 1181 102. **Afonine PV, Grosse-Kunstleve RW, Echols N, Headd JJ, Moriarty NW, Mustyakimov M,**
1182 **Terwilliger TC, Urzhumtsev A, Zwart PH, Adams PD.** 2012. Towards automated
1183 crystallographic structure refinement with phenix.refine. *Acta Crystallogr D Biol Crystallogr*
1184 **68**:352–367.
- 1185 103. **Emsley P, Lohkamp B, Scott WG, Cowtan K.** 2010. Features and development of Coot.
1186 *Acta Crystallogr D Biol Crystallogr* **66**:486–501.
- 1187 104. **Pei J, Kim B-H, Grishin NV.** 2008. PROMALS3D: a tool for multiple protein sequence and
1188 structure alignments. *Nucleic Acids Res* **36**:2295–2300.
- 1189 105. **Krissinel E, Henrick K.** 2004. Secondary-structure matching (SSM), a new tool for fast
1190 protein structure alignment in three dimensions. *Acta Crystallogr D Biol Crystallogr* **60**:2256–
1191 2268.
- 1192 106. **Hough MA, Wilson KS.** 2018. From crystal to structure with CCP4. *Acta Crystallogr D Struct*
1193 *Biol* **74**:67–67.
- 1194 107. **Robert X, Gouet P.** 2014. Deciphering key features in protein structures with the new
1195 ENDscript server. *Nucleic Acids Res* **42**:W320–4.
- 1196 108. **Sarkar S, Witham S, Zhang J, Zhenirovskyy M, Rocchia W, Alexov E.** 2013. DelPhi Web
1197 Server: A comprehensive online suite for electrostatic calculations of biological
1198 macromolecules and their complexes. *Commun Comput Phys* **13**:269–284.
- 1199 109. **Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, Ferrin TE.**
1200 2004. UCSF Chimera: a visualization system for exploratory research and analysis. *J*
1201 *Comput Chem* **25**:1605–1612.

Table 1

	wt IN A (NKNK)	IN A D116A	NQNQ	NQNK	NKNQ
Observed levels of integration 1 (relative to wt IN A) <i>Values from Figure 5</i>	1	0.00 ± 0.00	0.00 ± 0.00	0.03 ± 0.02	0.24 ± 0.09
Theoretical levels of 2LTRc 2 (relative to IN D116A) <i>see Materials and Methods</i>	0.20	1	1.00 ± 0.00	0.98 ± 0.02	0.81 ± 0.07
Observed levels of 2LTRc 3 (relative to IN D116A) <i>Values from Figure 7A</i>	0.20	1	0.33 ± 0.08	0.30 ± 0.09	0.28 ± 0.11
Efficiency of nuclear import 4 (relative to IN D116A) <i>Ratio values line 3 / values line 2</i>	1.00	1	0.33 ± 0.08	0.31 ± 0.09	0.35 ± 0.14
Ratio of PJ/2LTRc 5 (relative to IN D116A) <i>Values from figure 7B</i>	0.54 ± 0.02	1	0.93 ± 0.18	0.87 ± 0.19	0.66 ± 0.17
Decrease of PJ/2LTRc 6 (relative to IN D116A) <i>= 1-values in line 5</i>	0.46 ± 0.10	0	0.07 ± 0.01	0.13 ± 0.03	0.34 ± 0.09
Efficiency of 3' processing 7 (relative to wt IN A) <i>= values in line 6 / 0.46</i>	1.00 ± 0.22	0	0.15 ± 0.03	0.28 ± 0.06	0.74 ± 0.19
Expected levels of integration 8 (relative to wt IN A) <i>Product of values in lines 4 and 7</i>	1.00 ± 0.22	0	0.05 ± 0.01	0.09 ± 0.03	0.26 ± 0.11

In grey are given the values for the standard deviation (SD). For lines 1, 3 and 5 SD values are derived from the experimental values; for line 2, SD is given by 1-0.2 multiplied by the corresponding SD value from line 1; in line 4 SD is given by

$((SD_{line3}/average_{line3})^2 + (SD_{line2}/average_{line2})^2)^{1/2} \times average_{line4}$; in line 6 $SD=(SD_{line5}/average_{line5}) \times average_{line6}$; for line 7 $SD=SD_{line6}/0.46$; in line 8 SD is given by $((SD_{line7}/average_{line7})^2 + (SD_{line4}/average_{line4})^2)^{1/2} \times average_{line8}$.

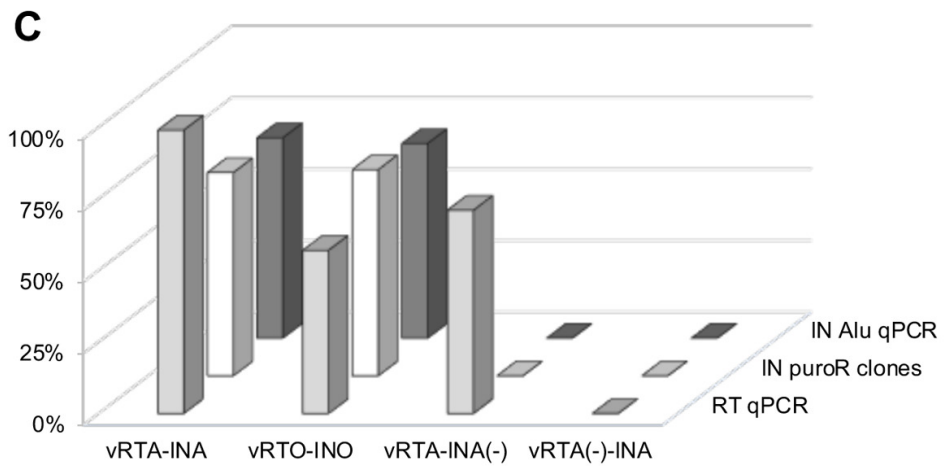
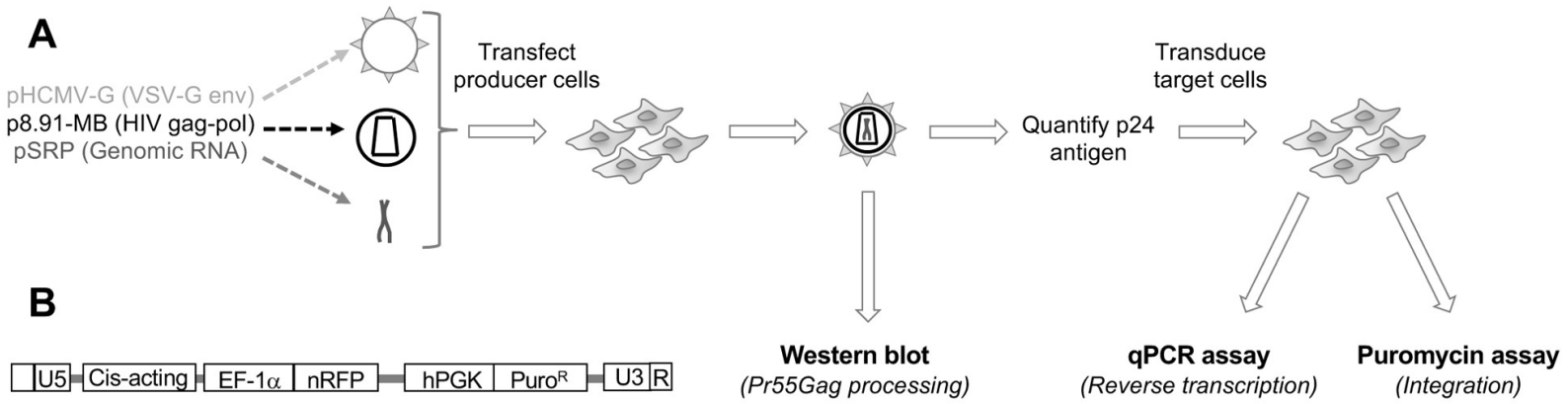


Figure 1

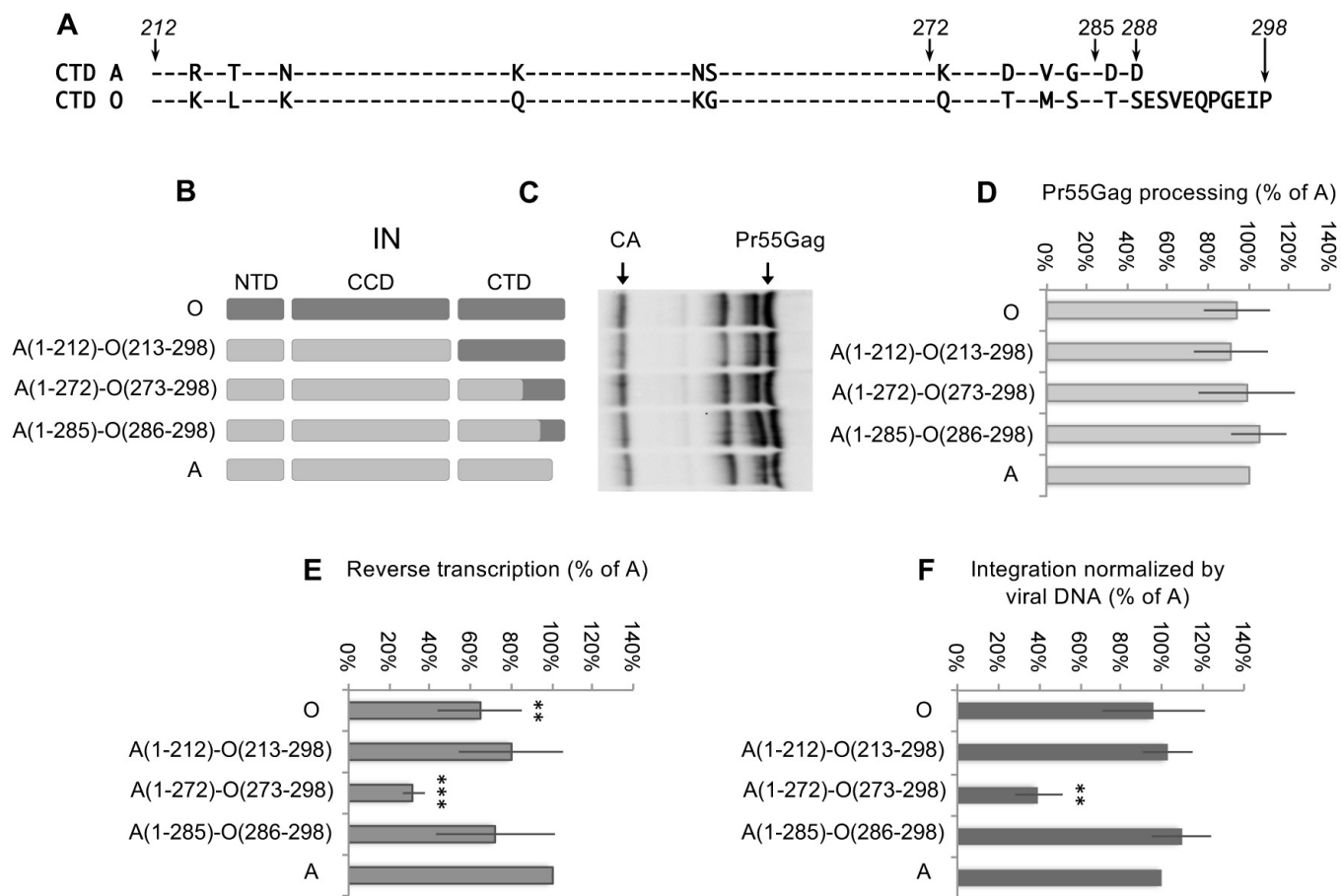


Figure 2

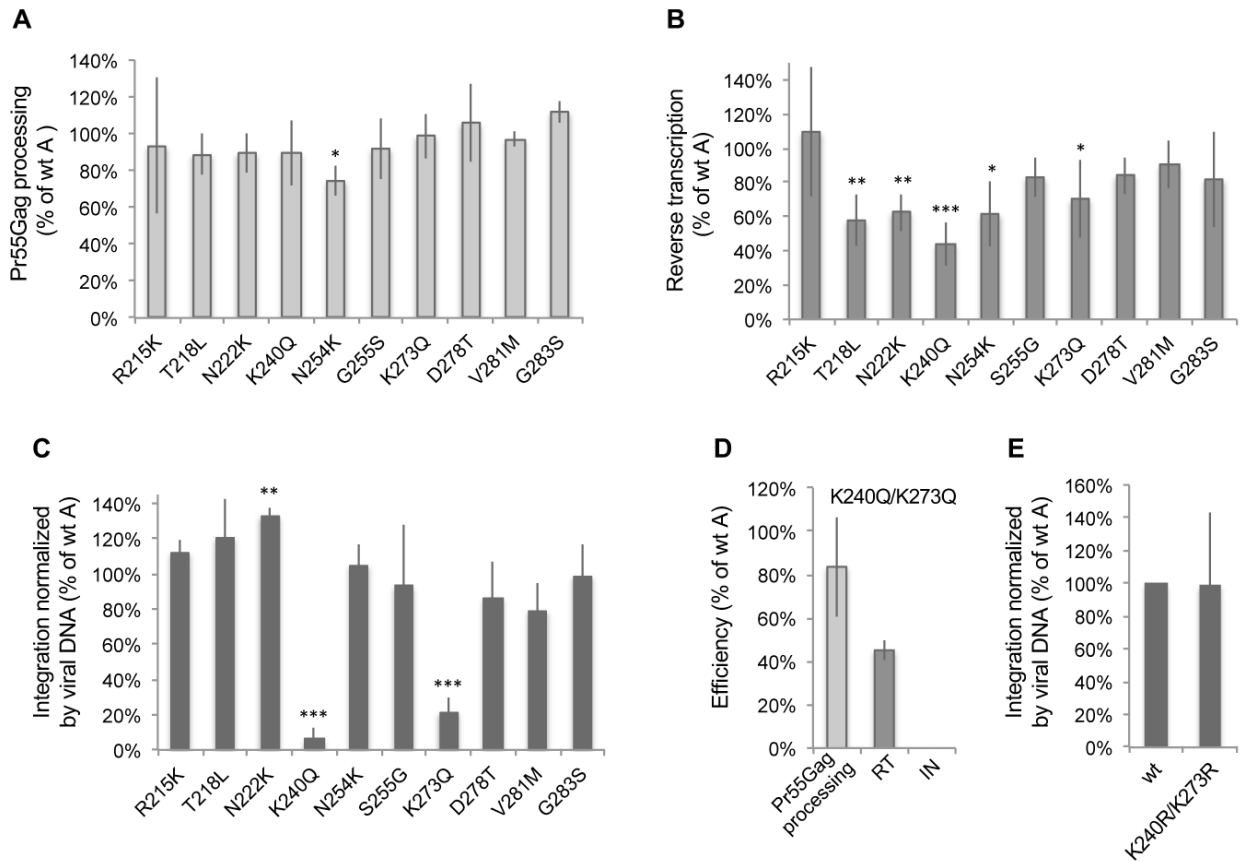


Figure 3

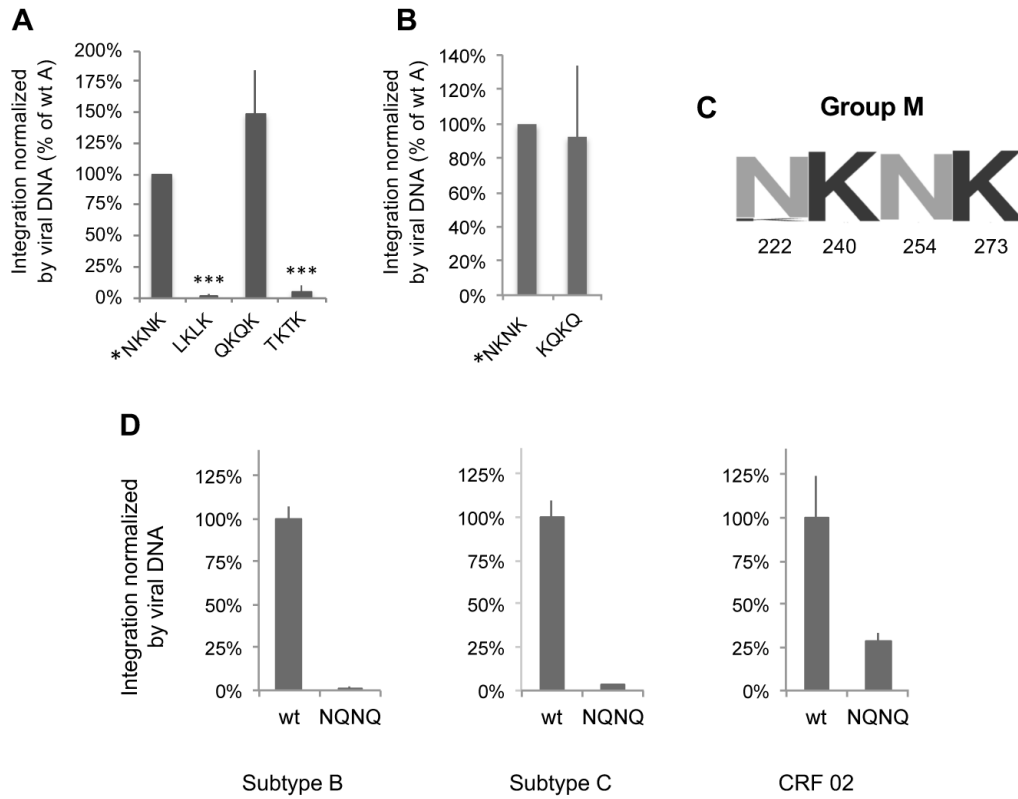


Figure 4

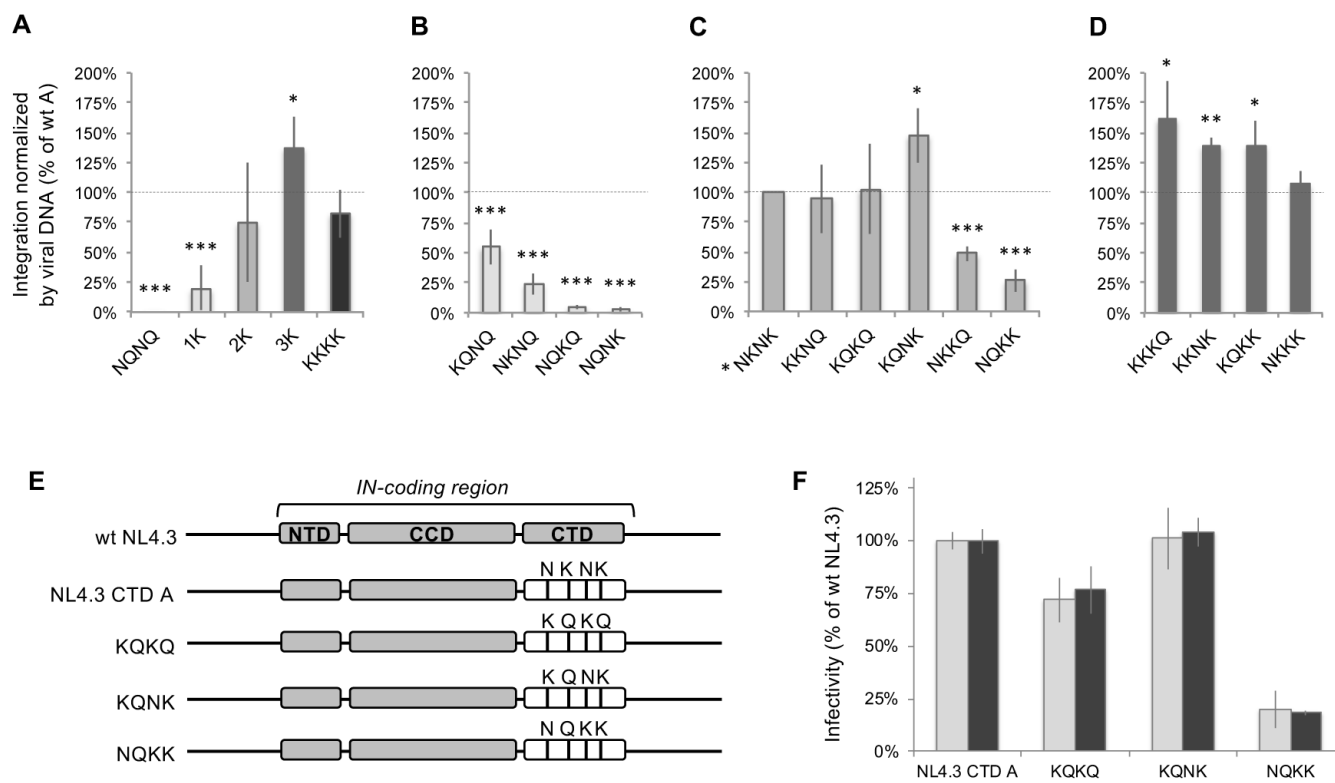


Figure 5

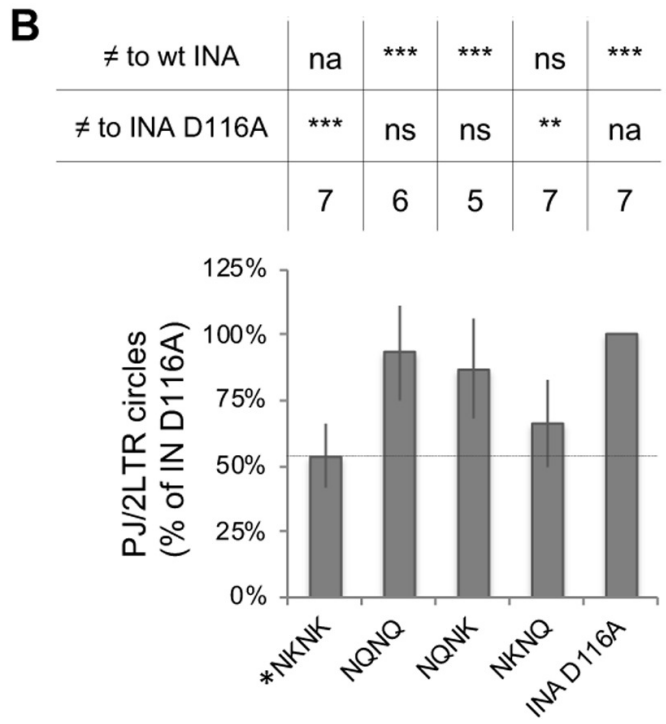
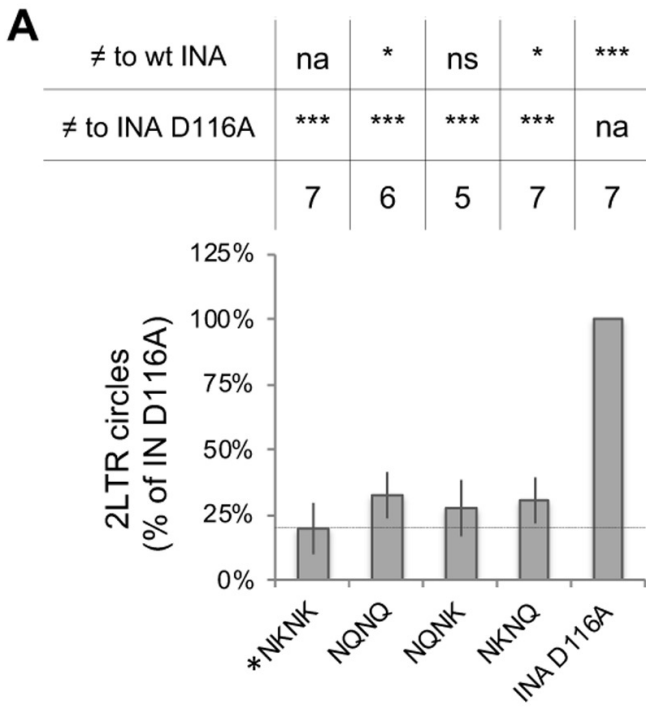


Figure 6

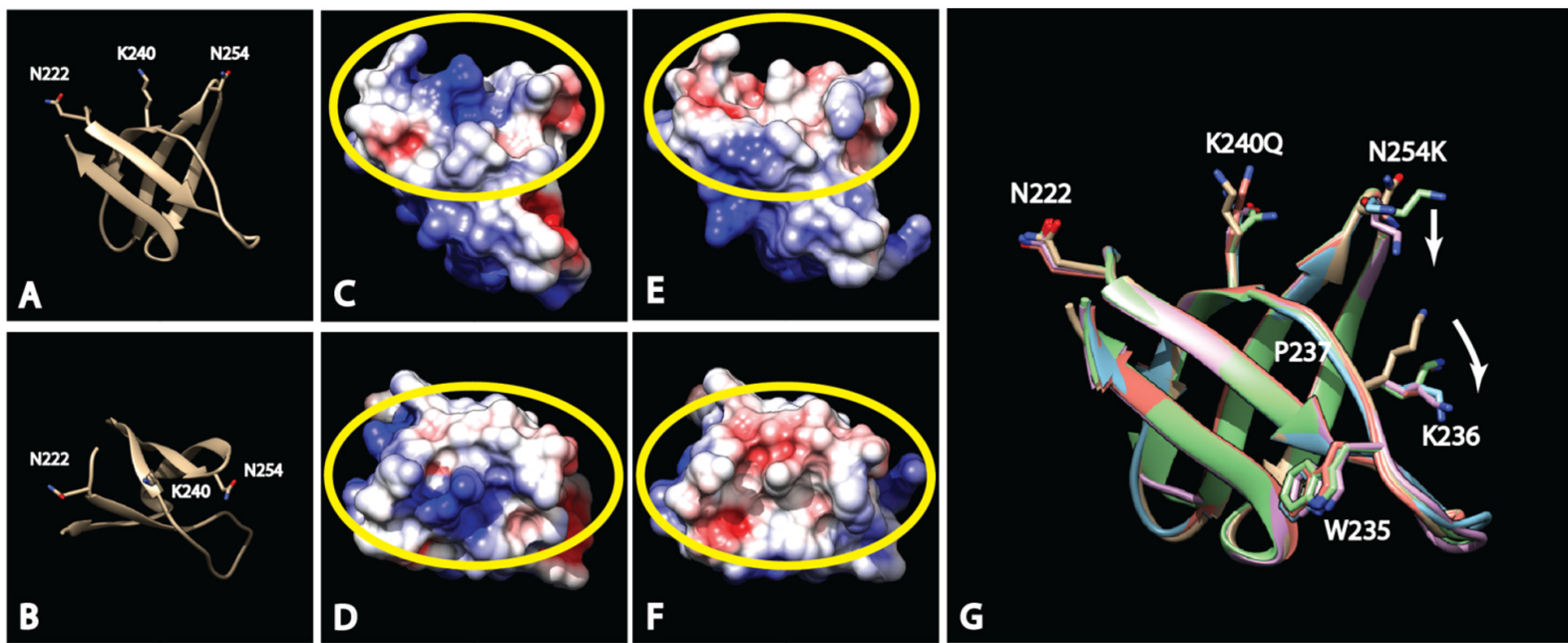


Figure 7

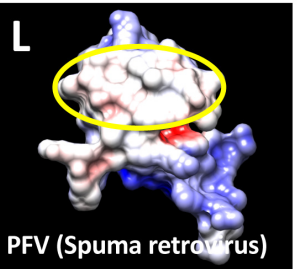
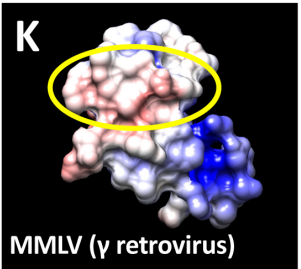
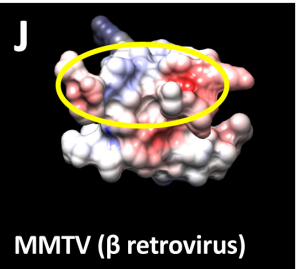
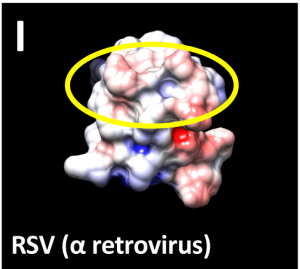
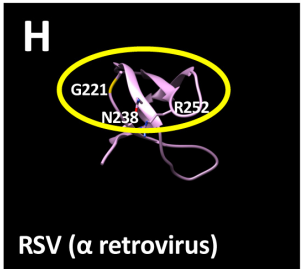
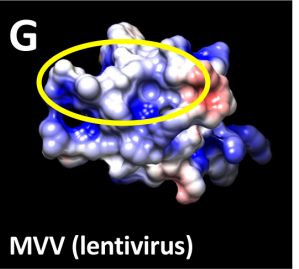
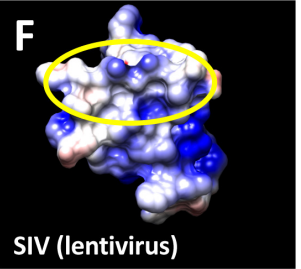
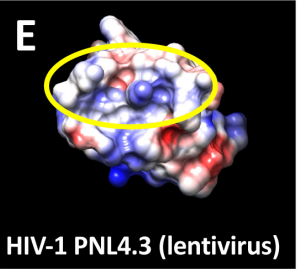
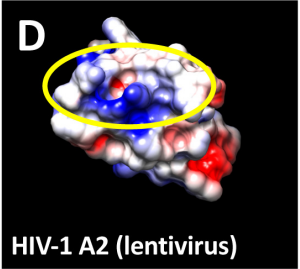
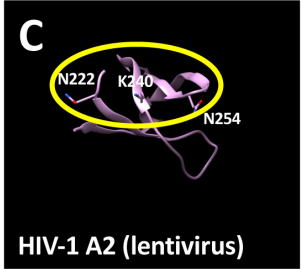
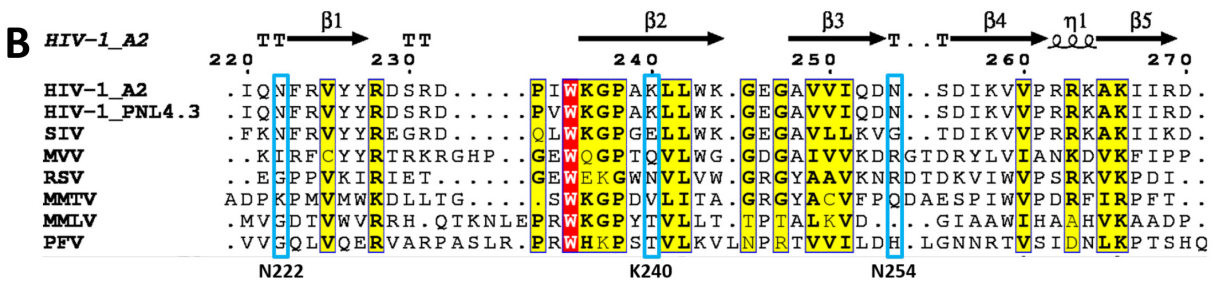
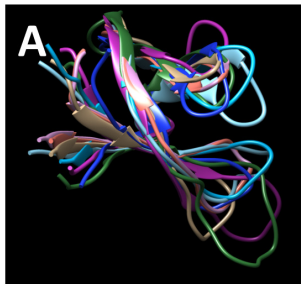


Figure 8