

Supporting Information

Lim et al. 10.1073/pnas.1001402107

SI Text

Supporting Materials and Methods. Mapping of viral integration sites.

A previously developed, high-throughput method for identifying human immune deficiency virus integration sites (1) was adapted to our system as follows. The genomic DNA of infected cells was isolated using Qiagen mini DNA kit (QIAGEN). The genomic DNA was then digested with various enzymes—BamH I, Hinc II, HpyCH4 III, HpyCH4 IV, Mse I, Hind III and Taq^{AI} [New England Biolabs (NEB)]—to minimize potential biases introduced by the choice of a few restriction enzymes. The fragmented genomic DNA was then used as the template for linear amplification with a biotinylated primer, bMLV, (5'-b-ATTTGTTAAA-GACAGGATATCAGTGGTCCAG-3') complementary to a sequence upstream of the 3' long terminal repeat (LTR) of integrated viral genome. The resulting amplification product was selectively concentrated using Dynabeads M-280 Streptavidin (Invitrogen) after PCR cleanup and then restricted by Mme I (NEB). The restricted fragments containing viral-host genome junctions were ligated to preannealed linker DNA (linkA, 5'-CGGATCCCGCATCAT-3', linkB, 5'-TGTCACACCTGGAGATATGATGCGGGATCCGNN-3') and then used for PCR with primers annealing to the 3' LTR and the linker (U5_5L, 5'-CTCCTGGATCCCCTCTTGCAGTTGCATCCGICTT-3', hi-linker, 5'-TCACACCTGGAGATATGATGCGG-3'). The amplified viral-host genome junctions were cloned into pBS SK SP plasmid after restriction with BamH I (NEB) (the BamH I sites within primers is underlined above) and then sequenced. As a complementary approach, the conventional method for cloning viral-host genome junctions was also applied here as described in detail elsewhere (2). The retroviral-host junction sequences were mapped to locations on the human genome (February 2009 assembly) using the BLASTN program on the Ensembl genome browser web site (www.ensembl.org). The chromosomal location mapping was considered for further statistical analysis only when the following conditions were satisfied: (i) no deletion was found before and at the CA sequence in the 3' end of 3' LTR (Fig. S1B) and (ii) a single chromosomal location can be assigned. Based on the chromosomal locations various genomic annotations for each retroviral integration site were made via genome browser web sites [Ensembl (www.ensembl.org) and UCSC genome browser (genome.ucsc.edu)].

1. Kim S, Kim Y, Liang T, Sinsheimer JS, Chow SA (2006) A high-throughput method for cloning and sequencing human immunodeficiency virus type 1 integration sites. *J Virol* 80:11313–11321.

Probability of having retroviral integrations at the same locations on the human genome. The probability that both 16.p12.zfd1 and 375.RT.zfd1 clones integrated into the same location on the human genome is calculated as follows:

The size of human genome = 3.1×10^9 .

For 16.p12.zfd1, the number of distinct virus-host junction sequences matching a single location on the human genome = 2

For 375.RT.zfd1, the number of distinct virus-host junction sequences matching a single location on the human genome = 4

Multiple virus-host junctions with the same sequence for each mutant case were counted one.

$$\begin{aligned} P &= {}_2C_1 \times (4/(3.1 \times 10^9)) \times ((3.1 \times 10^9 - 4)/(3.1 \times 10^9)) \\ &\sim {}_4C_1 \times (2/(3.1 \times 10^9)) \times ((3.1 \times 10^9 - 2)/(3.1 \times 10^9))^3 \\ &= 2.6 \times 10^{-9}. \end{aligned}$$

The probability that both 273.IN.zfd2 and 375.RT.zfd2 clones integrated into four common locations on the human genome is calculated as follows:

For 273.IN.zfd2, the number of distinct virus-host junction sequences matching a single location on the human genome = 22

For 375.RT.zfd2, the number of distinct virus-host junction sequences matching a single location on the human genome = 14

Multiple virus-host junctions with the same sequence for each mutant case were counted one.

$$\begin{aligned} P &\sim {}_{14}C_1 \times (22/(3.1 \times 10^9)) \times {}_{13}C_1 \\ &\times (21/(3.1 \times 10^9)) \times {}_{12}C_1 (20/(3.1 \times 10^9)) \\ &\times {}_{11}C_1 \times (19/(3.1 \times 10^9)) \times ((3.1 \times 10^9 - 22)/(3.1 \times 10^9))^{10} \\ &\sim {}_{22}C_1 \times (14/(3.1 \times 10^9)) \times {}_{21}C_1 \times (13/(3.1 \times 10^9)) \times {}_{20}C_1 \\ &\times (12/(3.1 \times 10^9)) \times {}_{19}C_1 \times (11/(3.1 \times 10^9)) \\ &\times ((3.1 \times 10^9 - 14)/(3.1 \times 10^9))^{18} = 4.6 \times 10^{-29}. \end{aligned}$$

2. Wu X, Li Y, Crise B, Burgess SM (2003) Transcription start regions in the human genome are favored targets for MLV integration. *Science* 300:1749–1751.

Table S2. Virus-host genome junction sequences identified for 273.IN.zfd2 and 273.IN.zfd3

Clones	Virus-host junction sequence	The number of junctions with the shown sequence	Total number of sequenced junctions for each mutant	% out of the total number of sequenced junctions	Integration position on the human genome	
					Chromosome	Position
273.IN.zfd2	<u>ACATT</u> TGGTAGCTGGGATGTTAG	22	197	11.2	X	153357666
273.IN.zfd2	<u>ACAT</u> GATGGGCAAAGTCACCC	10	197	5.1	1	94142416
273.IN.zfd2	<u>ACATT</u> TGGGGCTAAACAAATTT	10	197	5.1	X	113308851
273.IN.zfd2	<u>ACATT</u> ACAATTAATAAGTAT	8	197	4.1	2	13496337
273.IN.zfd2	<u>ACATT</u> TTCTTTCTAAAGTGCCT	6	197	3.0	15	93845542
273.IN.zfd2	<u>ACATT</u> TCATTAAGCGTAGGGC	6	197	3.0	13	74454561
273.IN.zfd2	<u>ACATT</u> TGGGCTTGGTAGGATTGG	5	197	2.5	7	86384864
273.IN.zfd2	<u>ACATT</u> TGGGGCTTAATATTTTT	4	197	2.0	16	52422344
273.IN.zfd2	<u>ACAT</u> CTAGACAGTTCAGAAAAA	3	197	1.5	2	146335293
273.IN.zfd2	<u>ACAG</u> TGGAGGGCACAGGATAT	3	197	1.5	21	41419380
273.IN.zfd2	<u>ACAT</u> CGGATCCCCGGGCTGCA	1	197	0.5	11	17549051
273.IN.zfd2	<u>ACATT</u> TGGGATCTGCCATATCA	1	197	0.5	11	116161720
273.IN.zfd2	<u>ACATT</u> TGGGGCTCTACATTACAG	1	197	0.5	16	58516469
273.IN.zfd2	<u>ACATT</u> TTCTTTCTAAGGTGCCT	1	197	0.5	4	173871541
273.IN.zfd2	<u>ACATT</u> TACTTTCTAAAGTGCCT	1	197	0.5	8	12828475
273.IN.zfd2	<u>ACATT</u> TGGGCTCCAAGATATTGC	1	197	0.5	7	62333487
273.IN.zfd2	<u>ACATT</u> TGGTAGCTGGTATGTTG	1	197	0.5	3	140396792
273.IN.zfd2	<u>ACATT</u> TGAAGTCCACAGCCTGCTGG	1	197	0.5	11	11519998
273.IN.zfd2	<u>ACATT</u> GAGTCAAACTAGAGCCT	1	197	0.5	15	88175538
273.IN.zfd2	<u>ACATT</u> TGGGGCCTGGACCACT	1	197	0.5	18	43408004
273.IN.zfd2	<u>ACATT</u> GGTAGCTGGGATGTTAGG	1	197	0.5	X	153357665
273.IN.zfd2	<u>ACATT</u> TGGAGACTAAATAAAAAC	1	197	0.5	1	97789945
273.IN.zfd3	<u>ACAT</u> TATTAATAGCAGTGTGCAC	18	101	17.8	5	93023433
273.IN.zfd3	<u>ACATT</u> ACCCTGGACCACTGATAT	15	101	14.9	14	52491298
273.IN.zfd3	<u>ACATT</u> GAGCCTGGACCACTGAT	7	101	6.9	16	75208997
273.IN.zfd3	<u>ACATT</u> GGAGGGGAGGAATGCCGT	5	101	5.0	18	46415241
273.IN.zfd3	<u>ACATT</u> TGGGGGCCAAGCTGCCT	5	101	5.0	16	56729532
273.IN.zfd3	<u>ACATT</u> ACCCTGGACCACTGATGT	3	101	3.0	8	143413073
273.IN.zfd3	<u>ACATT</u> TGGCGCCCCGAGTGAGGG	1	101	1.0	1	180199185
273.IN.zfd3	<u>ACATT</u> ACCCTGGACCACTGATATC	1	101	1.0	17	8050843
273.IN.zfd3	<u>ACATT</u> CCCTACCTACATTGTT	1	101	1.0	19	24183658

Virus-host genome junction sequences are shown, along with the fraction of total junctions corresponding to that integration position for each mutant. The underlined sequences are viral, and the others are from the host genome.