

Transmitted HIV-1 is more pathogenic in heterosexual individuals than homosexual men

Ananthu James^a and Narendra M. Dixit^{a,b,*}

^aDepartment of Chemical Engineering, Indian Institute of Science, Bengaluru, India

^bCentre for Biosystems Science and Engineering, Indian Institute of Science, Bengaluru, India

*Corresponding author; E-mail: narendra@iisc.ac.in

Manuscript details:

Abstract: 150 words

Text: ~1800 words

Figures: 3

Table: 1

References: 27

Methods: ~1000 words

Supplementary Materials:

Text notes: 2

Figures: 1

Tables: 7

References: 41

1 **Abstract**

2 Transmission bottlenecks introduce selection pressures on HIV-1 that vary substantially with the mode
3 of transmission. Recent studies on small cohorts have suggested that stronger selection pressures lead
4 to fitter transmitted/founder (T/F) strains. Manifestations of this selection bias at the population
5 level have remained elusive. Here, we analysed early CD4 cell count measurements reported from
6 ~340,000 infected heterosexual individuals (HSX) and men-who-have-sex-with-men (MSM), across
7 geographies, ethnicities and calendar years and found them to be consistently lower in HSX than
8 MSM ($P < 0.05$). The corresponding average reduction in CD4 counts relative to healthy adults was
9 86.5% in HSX and 67.8% in MSM ($P < 10^{-4}$). This difference could not be attributed to differences in
10 age, HIV-1 subtype, viral load, gender, ethnicity, time of transmission, or diagnosis delay across the
11 groups. We concluded that the different selection pressures arising from the different predominant
12 transmission modes have resulted in more pathogenic T/F strains in HSX than MSM.

13 Introduction

14 The bottlenecks in HIV-1 transmission result in a ‘selection bias’ favoring fitter transmitted/founder
15 (T/F) viruses over less fit ones^{1,2}. Several recent studies have presented evidence of genetic, pheno-
16 typic, and clinical manifestations of the selection bias in small cohorts^{1,3–6}. From 137 heterosexual
17 (HSX) donor-recipient pairs, T/F viruses were found to carry higher than average frequencies of
18 amino acids associated with high *in vivo* fitness¹. Similarly, from 127 discordant couples, lower vi-
19 ral replication capacity (vRC), indicative of lower viral fitness, early in infection was associated with
20 slower decline of CD4 T cell counts^{4,6}. The selection bias varies with the mode of transmission³.
21 The stronger the bottlenecks, the fitter the corresponding T/F viruses are likely to be^{1,2}. Anal inter-
22 course is over 10-fold more permissive on average than penile-vaginal intercourse⁷. Analysis of T/F
23 genomes from 131 subjects revealed that the T/F genomes were under greater positive selection in
24 heterosexual individuals (HSX), in whom the penile-vaginal mode predominates⁸, than homosexual
25 men, or men-who-have-sex-with-men (MSM), who transmit predominantly through anal intercourse³.
26 Among HSX, men had T/F viruses with higher predicted fitness *in vivo* than women¹, consistent with
27 the asymmetry of the bottlenecks between insertive and receptive penile-vaginal intercourse⁷.

28 An important question that follows is whether the differential selection bias across modes of trans-
29 mission is manifested at the wider population level. Such differential bias could contribute to varia-
30 tions in disease progression and treatment outcomes and underlie the diverse trajectories of the HIV-1
31 pandemic across infected groups in which different modes of transmission predominate.

32 Results and Discussion

33 To answer this question, we decided to compare early CD4 T cell count measurements between HSX
34 and MSM. Immediately following infection, CD4 T cell counts fall steeply, recover partially, and then
35 settle within a few weeks/months to a value smaller than in the pre-infection state⁹ (Fig. 1(a)).
36 Subsequent changes in the CD4 counts occur slowly, over many months to years. Thus, CD4 count
37 measurements made early in infection tend to be close to the value to which the counts settle after
38 the initial dynamics. These early CD4 counts are expected to be minimally affected by host-specific
39 adaptive mutations¹ and, therefore, representative of the fitness of the T/F strain in the recipient. The
40 fitter the strain, the lower would be the CD4 count. The CD4 count is also a more robust marker of
41 disease state than other commonly used markers such as set-point viral load (SPVL). High vRC of the

42 T/F viruses was associated with low CD4 counts at 3 months post-infection (which roughly coincides
43 with seroconversion) and rapid CD4 count decline for ~ 5 years, independently of SPVL^{4,6}.

44 HSX and MSM are the two major groups driving the global HIV-1 epidemic⁹. They use predominant
45 modes of transmission with a substantial difference in the selection bias⁷. Importantly, they display
46 little inter-mixing in most geographical regions. We inferred the latter from the distinct prevalence
47 of HIV-1 subtypes in the two groups, which we found across geographical regions and calendar years
48 (Fig. 2; Text S1; Tables S1 and S2). Together, these characteristics allow for the difference in the
49 selection bias to be sustained long-term, potentially amplified, and manifested in sample sizes large
50 enough for detection with statistical significance. We thus hypothesized that the stronger selection
51 bias associated with penile-vaginal transmission than anal transmission would result in lower early
52 CD4 counts in HSX than in MSM.

53 To test this hypothesis, we collated available data of CD4 count measurements either at serocon-
54 version or at diagnosis from all large studies, which amounted to a total of $\sim 340,000$ patients across
55 four geographical regions followed over a total period of nearly four decades, and examined the dif-
56 ferences between HSX and MSM (Methods; Table 1). We found that HSX consistently had lower CD4
57 counts than MSM (Fig. 1(b); Tables 1 and Tables S3-S5). For instance, measurements from $\sim 120,000$
58 patients across 21 countries in the European Union and European Economic Area (EU/EEA) indicated,
59 following population-weighted averaging of yearly data during 2010-2018, that the mean CD4 count
60 in MSM at diagnosis was ~ 440 cells/ μL , whereas it was substantially lower, ~ 300 cells/ μL , in HSX
61 ($P < 10^{-4}$)¹⁰. The numbers were similar in the preceding 5 year period (2002-2007) reported by a
62 smaller study involving a few thousand patients¹¹. In the UK, measurements from close to 9000 pa-
63 tients during 1990-1998 showed that the counts at diagnosis were ~ 330 cells/ μL in MSM and ~ 230
64 cells/ μL in HSX ($P < 10^{-3}$)¹². In China, during 2006-2012, the mean CD4 counts at diagnosis from
65 $\sim 180,000$ patients were ~ 370 cells/ μL in MSM and ~ 270 cells/ μL in HSX ($P < 10^{-4}$)¹³. Similarly, in
66 the US, from over 25,000 patients during 2006-2015, the counts at diagnosis were ~ 400 cells/ μL in
67 MSM and ~ 300 cells/ μL in HSX ($P < 10^{-4}$)¹⁴. We also examined/estimated the counts at seroconver-
68 sion where available. In the CASCADE study, involving ~ 4000 patients during 1979-2000 in Europe
69 and Australia, the mean cell counts at seroconversion were ~ 620 cells/ μL in MSM and ~ 590 cells/ μL
70 in HSX ($P = 0.027$)¹⁵. Further, using the reported diagnosis delays and the slopes of CD4 count decline

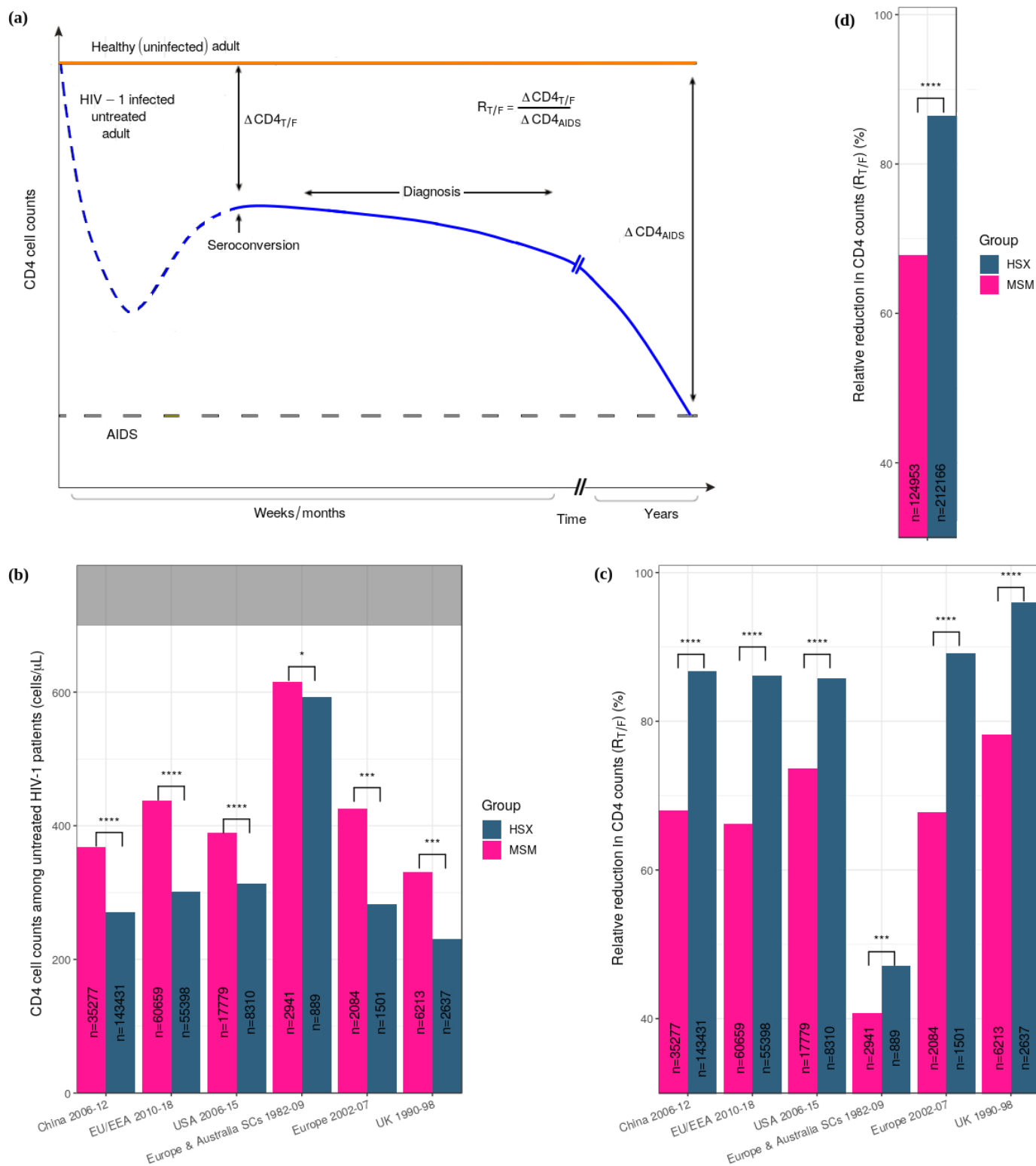


Figure 1: Early CD4 T cell counts and the associated relative reduction ($R_{T/F}$) in MSM and HSX. (a) Schematic of typical CD4 count changes post HIV-1 infection (blue), before (dashed) and after (solid) diagnosis/seroconversion. The reduction at diagnosis/seroconversion relative to uninfected individuals (orange) and that associated with AIDS (grey dashed line) yields $R_{T/F}$, the reduction attributable to the T/F virus. **(b)** Early mean CD4 cell counts and **(c)** the corresponding $R_{T/F}$ in untreated infected adult HSX and MSM from different geographical regions and calendar years (see Methods, Tables 1 and S3-S5 for details). The grey region indicates counts in uninfected, healthy individuals. **(d)** Population-weighted average of $R_{T/F}$ across all the datasets in (c). The sample sizes (n) are indicated. SCs indicate seroconverters. ****, ***, ** and * indicate $P < 10^{-4}$, $P < 10^{-3}$, $P < 10^{-2}$ and $P < 0.05$, respectively.

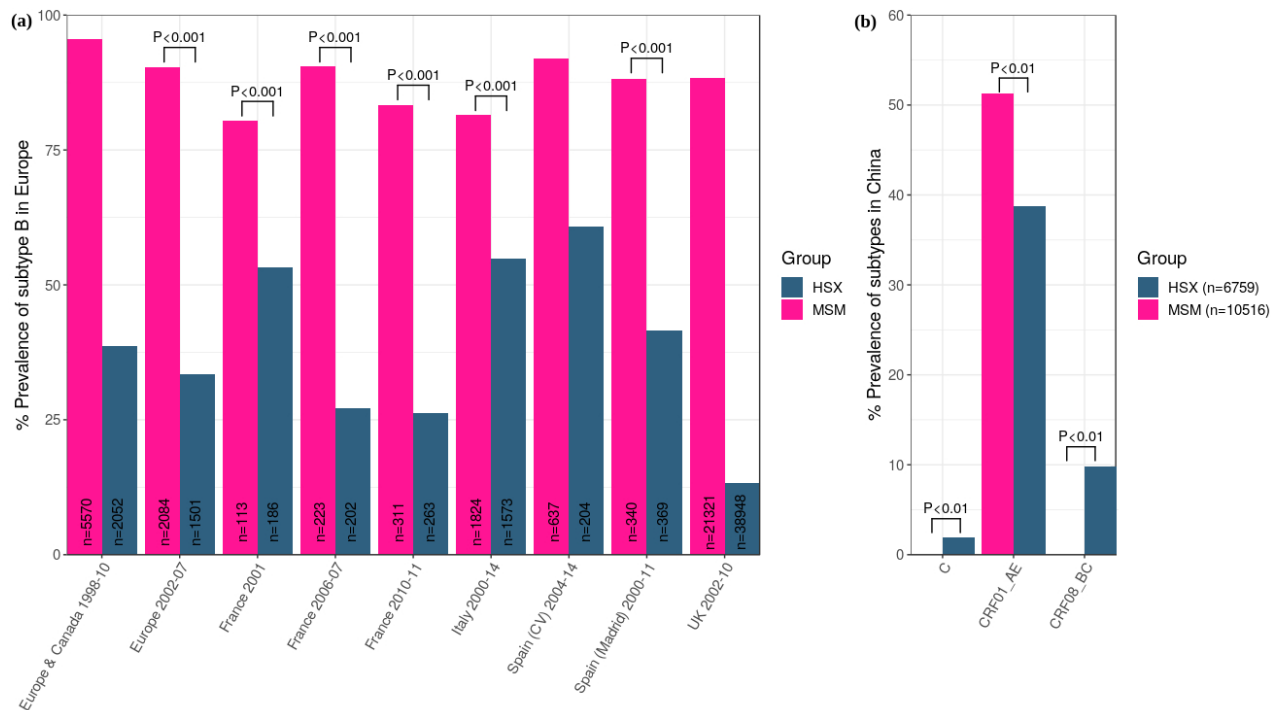


Figure 2: Subtype prevalence of HIV-1 in MSM and HSX populations. Prevalence of (a) subtype B in different regions in Europe and Canada and (b) all the subtypes in China. The sample sizes (n) along with the time periods of the surveys are indicated. P values are listed where available in the original sources. Sources of the data and additional details are in Tables S1 and S2. The different prevalence indicates little mixing between MSM and HSX in the populations studied.

71 in the US population above¹⁴, we estimated that the cell counts at seroconversion, for the age group
 72 13-29 years, were ~ 550 cells/ μ L in MSM and ~ 480 cells/ μ L in HSX ($P < 10^{-4}$) (Methods). Remark-
 73 ably, we did not find any large study (sample size $\gtrsim 1000$) that reported higher early CD4 cell counts
 74 in HSX than MSM.

75 While the evidence from absolute CD4 count comparisons was thus overwhelming, differences in
 76 CD4 counts in healthy (uninfected) individuals across gender, ethnicity and geographical regions
 77 could render absolute CD4 counts only an approximate measure of the fitness of the T/F strains. Two
 78 individuals may have similar early CD4 counts but may still have been infected by T/F strains of
 79 different fitness if their pre-infection CD4 counts were different, with the individual with the higher
 80 pre-infection count infected by the fitter T/F strain. To overcome this limitation, we constructed a
 81 metric to quantify the relative reduction in the CD4 cell count, R , corresponding to the absolute CD4
 82 count T as $R = \frac{T_{healthy} - T}{T_{healthy} - T_{AIDS}} \times 100 = \frac{\Delta CD4}{\Delta CD4_{AIDS}} \times 100$, where $T_{healthy}$ was the count pre-infection, and
 83 $T_{AIDS} = 200$ cells/ μ L the count defining AIDS. Thus, R was 0% when $T = T_{healthy}$ and 100% when
 84 $T = T_{AIDS}$ and decreased linearly with T between these extremes. Choosing $T_{healthy}$ specific to the

Table 1. Early CD4 cell counts in infected adults at diagnosis or seroconversion, estimated from data collated from several studies (see text and Tables S3-S5). The corresponding relative reduction in CD4 count, $R_{T/F}$, calculated as described in the text, are also listed. P values for comparisons of the CD4 counts (and $R_{T/F}$) between MSM and HSX indicate a significantly higher CD4 counts and lower $R_{T/F}$ in MSM throughout. The last row represents the population-weighted average of all the datasets.

Sample & duration	Risk group	CD4 counts (<i>cells/μL</i>)		Sample size (<i>n</i>)	$R_{T/F}$ (%)	
		Mean (SD)	P value		Mean (SD)	P value
EU/EEA 2010 – 18	MSM*	437 (397)	0.00 ^{††}	60,659	66.2 (58.7)	0.00 ^{††}
	HSX	302 (464) ^d		55,398	86.2** (63.0)	
	HSX men*	270 (464) ^d	2.0 × 10 ⁻⁸¹	27,822	90.0 (66.3)	1.7 × 10 ⁻⁵¹
	HSX women	335 (325)		27,576	82.9 (42.0)	
USA 2006 – 15	MSM (6 m.o. diagnosis)	390 (243)	3.0 × 10 ⁻⁹⁷	17,779	73.7 (34.4)	5.2 × 10 ⁻¹⁴³
	HSX (6 m.o. diagnosis)	314 (284)		8,310	85.8 (35.9)	
	MSM (SCs, age: 13-29)	553 (315)	1.2 × 10 ⁻¹⁸	6,328	51.0 (45.5)	1.8 × 10 ⁻³⁷
	HSX (SCs, age: 13-29)	483 (377)		2,958	64.7 (49.0)	
China 2006 – 12	MSM	368 (222)	0.00 ^{††}	35,277	68.0 (45.1)	0.00 ^{††}
	HSX	270 (260)		143,431	86.7 (50.0)	
Europe 2002 – 07	MSM	426 (242)	$c < 10^{-3}$	2,084	67.8 (37.5)	2.5 × 10 ⁻⁶³
	HSX	283 (258)		1,501	88.8 (35.3)	
UK 1990 – 98	MSM	331 (242) ^d	$c < 10^{-3}$	6,213	78.2 (42.0)	3.8 × 10 ⁻⁹²
	HSX	230 (258) ^d		2,637	96.0 (34.7)	
Europe & Australia (CASCADE) SCs 1979 – 00	MSM [†] (age < 40)	621 (323)	4.3 × 10 ⁻³	2,570	40.0 (53.7)	0.010
	HSX men ^a (age < 40)	576 (327)		428	46.4 (52.7)	
	HSX women ^{†,a} (age < 40)	623 (229)		349	46.4 (38.7)	
	MSM [‡] (age > 40)	578 (221)	0.017	371	46.2 (40.1)	0.073
	HSX men ^b (age > 40)	534 (134)		62	52.4 (29.1)	
	HSX women ^{‡,b} (age > 40)	580 (113)		50	51.8 (27.1)	
	MSM (total)	616 (323) ^d	0.027	2,941	40.7 (53.7)	6.9 × 10 ⁻⁴
	HSX (total)	592 (327) ^d		889	47.1 (51.7)	
Overall	MSM	–		124,953	67.8 (58.7) ^d	0.00 ^{††}
	HSX	–		212,166	86.5 (63.0) ^d	

* $P = 0.00^{\dagger\dagger}$ for the comparisons of both the cell counts and $R_{T/F}$ between MSM and HSX men.

** $R_{T/F}$ in the EU/EEA region changed from 86.2% (SD = 63.0%) to 84.9% (SD = 69.5%) in HSX when we used CD4 counts in healthy individuals from Tanzania (representative of sub-Saharan Africa, 746 cells/μL (Table S6)), instead of 941 cells/μL from Italy (representative of EU/EEA), corresponding to the 35% of HSX, still substantially higher than $R_{T/F}$ of 66.2% in MSM ($P = 0.00^{\dagger\dagger}$).

^{††}These P values were below the lower representation limit ($= 2.23 \times 10^{-308}$) of R (and Excel).

[†] $P = 0.44$ for CD4 count and $P = 3.1 \times 10^{-3}$ for $R_{T/F}$ comparisons between MSM and HSX women aged < 40 years.

[‡] $P = 0.46$ for CD4 count and $P = 0.10$ for $R_{T/F}$ comparisons between MSM and HSX women aged > 40 years.

^a $P = 0.010$ for CD4 count and $P = 0.50$ for $R_{T/F}$ comparisons between HSX men and HSX women aged < 40 years.

^b $P = 0.026$ for CD4 count and $P = 0.46$ for $R_{T/F}$ comparisons between HSX men and HSX women aged > 40 years.

^cReported in the original sources.

^dSee Methods and the footnotes in Tables S3 and S5.

85 respective geographies, ethnicities, and genders (Table S6), we estimated R corresponding to the
86 early cell count measurements above, which we denoted as $R_{T/F}$, indicative of the relative reduction
87 in CD4 count due to the T/F virus (Fig. 1(c)). The higher the $R_{T/F}$, the fitter would be the T/F
88 strain, regardless of the pre-infection CD4 count, rendering $R_{T/F}$ a more robust marker of T/F viral
89 fitness than the associated early absolute CD4 counts. (Note that $R_{T/F}$ is a static measure and is not
90 indicative of the ‘speed’ of disease progression; cell count decline can be faster despite higher early
91 CD4 counts in MSM than HSX^{15,16}.)

92 We found that in EU/EAA, during 2010-18, $R_{T/F}$ was 86.2% in HSX and 66.2% in MSM ($P < 10^{-4}$).
93 During 2002-07, these numbers were 88.8% and 67.8% ($P < 10^{-4}$), respectively. The corresponding
94 numbers were 96.0% and 78.2% in the UK ($P < 10^{-4}$), and 86.7% and 68.0% in China ($P < 10^{-4}$).
95 In the US, the difference was smaller but still substantial, with $R_{T/F}$ of 85.8% in HSX and 73.7% in
96 MSM ($P < 10^{-4}$). At seroconversion, these numbers were 64.7% and 51.0%, respectively ($P < 10^{-4}$).
97 For the seroconverters from the CASCADE study, the trend was consistent, with $R_{T/F}$ of 47.1% in HSX
98 and 40.7% in MSM ($P < 10^{-3}$). Overall, thus, $R_{T/F}$ comparisons showed more significant differences
99 between MSM and HSX than absolute CD4 count comparisons (Fig. 1(b) and (c)). Further, $R_{T/F}$ al-
100 lowed comparison across the different datasets. Thus, while the HSX all had $R_{T/F} > 85\%$ at diagnosis,
101 the MSM displayed a range from $\sim 65\%$ to a little under 80%. We could also combine the datasets,
102 including those at diagnosis and seroconversion, and estimate an overall $R_{T/F}$. Using a population-
103 weighted average across the datasets, we estimated the overall $R_{T/F}$ to be 86.5% in HSX and 67.8%
104 in MSM ($P < 10^{-4}$) (Fig. 1(d)). This overall comparison provides strong evidence of greater cell count
105 reduction due to, and hence greater pathogenicity of, the T/F viruses in HSX than in MSM.

106 To attribute the differences in $R_{T/F}$ between HSX and MSM to the differential selection bias at
107 transmission in the two groups, we considered and ruled out all the major potential confounding
108 factors. First, MSM are typically diagnosed at a younger age than HSX. In the two European studies,
109 MSM were 5 (Table S5)¹⁰ and 1.6 years¹¹ younger on average than HSX at diagnosis. Given the cell
110 count decrease of ~ 7 cells/ μ L per year of age at diagnosis¹², the CD4 counts should have been higher
111 in MSM by only ~ 35 and ~ 11 cells/ μ L, whereas they were higher by 135 and 143 cells/ μ L (Fig. 1(b)),
112 respectively, a difference that could not be explained by the age at diagnosis. Second, MSM are often
113 predominantly infected by subtype B¹⁷, whereas HSX are by subtypes B and C (Fig. 2; Text S1). This

114 subtype difference should have resulted in lower CD4 counts in MSM than HSX because of the higher
115 virulence of subtype B^{18,19}, a trend opposite of what is observed. Moreover, in the US where subtype
116 B dominates both HSX and MSM (Text S1), $R_{T/F}$ was lower among MSM (Fig. 1(c)). In agreement,
117 an independent study found that subtype B T/F viruses had higher fitness among HSX than MSM³.
118 Third, the CD4 counts could not be explained as an indirect manifestation of variations in SPVL; in
119 the European study, CD4 counts were higher in MSM despite higher SPVL in MSM than HSX (Table
120 S3). Fourth, healthy men had lower CD4 counts than HSX and healthy women everywhere except
121 China (Table S6), and infected HSX men displayed higher $R_{T/F}$ than MSM (Table 1 and Fig. S1),
122 two reasons to rule out gender as the cause of lower $R_{T/F}$ in MSM. Fifth, in Europe (EU/EEA), while
123 MSM are predominantly Caucasian, 30-35% of infected HSX are of sub-Saharan African origin^{10,11}.
124 In China, however, where no differences in ethnicity exist between MSM and HSX, a substantial
125 difference in $R_{T/F}$ is seen between them (Fig. 1(c)), ruling out ethnicity as a confounding factor.
126 Further, accounting for baseline CD4 count differences across ethnicities in EU/EEA did not alter our
127 findings (Table 1). Sixth, early onward transmission may limit donor-specific adaptations in the T/F
128 strain and allow it to cause more severe cell count reduction in the recipient. Early transmissions,
129 however, are more common to MSM than HSX^{18,20}, in keeping with the greater association of MSM
130 with transmission clusters¹⁷ (Fig. 3; Table S7), and should have led to higher $R_{T/F}$ in MSM than HSX,
131 in contrast to our findings. Seventh, although MSM tend to be diagnosed earlier than HSX¹⁴ and
132 may thus suffer a lower loss of CD4 counts at diagnosis, the differences are seen also in CD4 counts
133 at seroconversion^{14,15}, which would occur at similar times post infection in the two groups. Besides,
134 MSM had lower cell counts in China too, where, owing to social stigma, MSM may not get diagnosed
135 earlier than HSX¹³. The difference in $R_{T/F}$ between MSM and HSX was thus not attributable to any
136 of the above factors. We concluded therefore that the difference originated from the variations in the
137 fitness of the T/F strains in the two groups arising from the different selection biases at transmission.

138 Our findings establish the selection bias at transmission as an important underlying factor shaping
139 HIV-1 adaptation at the population level. The differential adaptation of HIV-1 to MSM and HSX,
140 which in most geographical regions show little inter-mixing, may have led over the years to the
141 selection and, possibly, fixation of different adaptive mutations in the T/F viruses in the two groups.
142 Genetic differences have been observed between T/F strains in MSM and HSX in small cohorts³.

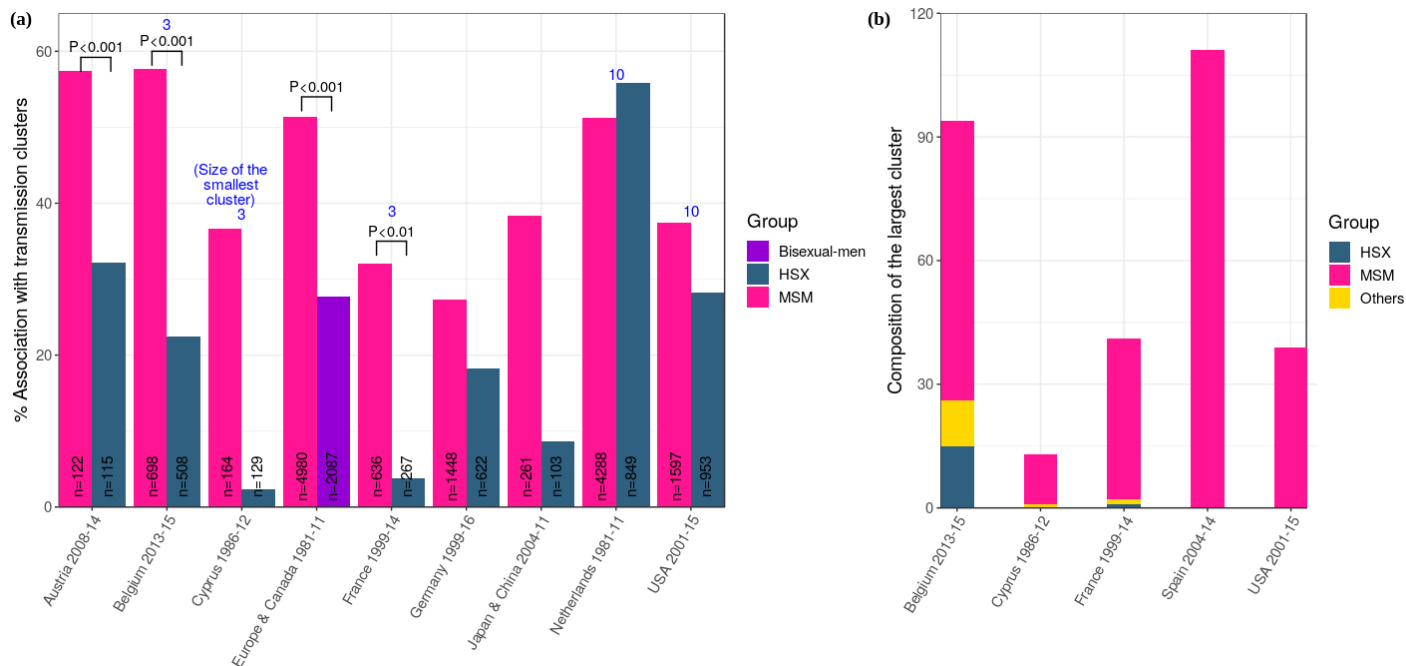


Figure 3: Association with transmission clusters. (a) The fraction of HSX and MSM (or bisexuals in one case) associated with transmission clusters and (b) the composition of the largest clusters in different geographical regions. The sample sizes (n) along with the time periods of the surveys are indicated. P values and the minimum sizes of the clusters (blue text) where available from the original sources are listed. Sources of the data and additional details are in Table S7. MSM are thus far more likely to be associated with clusters and also tend to form large clusters.

143 Future studies may establish them at the population level, as sequencing technologies that allow facile
 144 identification of T/F viruses emerge. The technologies may also serve to elucidate such differences
 145 between other infected groups, which are likely to be present to lower degrees than between MSM
 146 and HSX, depending on the differences in the selection bias between the groups, the exclusivity of
 147 the associated modes of transmission, and the extent of mixing between the groups. Our findings
 148 also suggest that heritable viral traits such as SPVL²¹ may have evolved differently in MSM and HSX,
 149 potentially driving differential spread of the HIV-1 epidemic in the two groups. The extent of these
 150 differences may determine whether intervention strategies, including the development and use of
 151 preventive vaccines, may have to be tailored to individual infected groups.

152 Methods

153 **Data of CD4 counts.** To test our hypothesis that early CD4 counts in HSX would be higher than
 154 in MSM at the population level, we collated data from all large studies ($n \gtrsim 1,000$) that reported
 155 CD4 counts either at diagnosis or seroconversion in both these groups. The data are summa-
 156 rized along with our analysis in Table 1 and details are in Tables S3-S5. From reports on coun-

157 tries in the EU/EEA and China^{10,13}, we digitized the median CD4 counts using WebPlotDigitizer
158 (<https://automeris.io/WebPlotDigitizer>). For our analysis, we averaged the data over the study dura-
159 tion. To obtain sample sizes, we multiplied the diagnosed cases with the reported fraction of diagnoses
160 contributing to the annual CD4 counts in the entire EU/EEA (Table S5). The fraction was assumed to
161 be the same across the risk groups and the set of 21 countries studied. We also assumed the propor-
162 tions of men and women in HSX to remain the same during 2010-18. In the CASCADE study¹⁵, which
163 segregated data into age groups, we averaged over age groups. To obtain the population-weighted
164 average CD4 counts, we assumed that the proportions of the populations in the different transmis-
165 sion categories were the same across age groups and that the fractions of men and women remained
166 conserved (except in MSM and hemophiliacs) (Table S3). To calculate $R_{T/F}$, we also collated data
167 of CD4 counts from healthy, uninfected adults in the USA, UK, Italy (which was used for the three
168 studies involving European populations), Tanzania, and China, which are listed in Table S6. For $R_{T/F}$
169 calculations pertaining to the UK, CD4 counts from healthy MSM and HSX were available, which we
170 used. We found the counts in MSM comparable to those from healthy HSX men. As a result, for other
171 populations, we used the cell counts for healthy HSX men where counts from healthy MSM were
172 unavailable.

173 **Estimation of mean CD4 counts and their standard deviations.** When the median, m , and in-
174 terquartile range (IQR), (q_l, q_u) , of CD4 counts were available, we estimated the corresponding
175 mean, μ , and standard deviation (SD), σ , using $\mu = \frac{m+q_u+q_l}{3}$ and $\sigma = \frac{q_u-q_l}{1.35}$, following the widely
176 used method²² applicable to large sample sizes, as considered here. When 95% confidence intervals
177 (CIs), (c_l, c_u) , were available instead of IQR, we evaluated SD using another method²³ which yielded
178 $\sigma = \frac{\sqrt{n}(c_u-c_l)}{3.92}$ when the sample size $n \gtrsim 100$. When IQR was unavailable, we approximated the medi-
179 ans as the means, assuming the distributions to be normal. For data from China and EU/EEA, where
180 σ was available for the total population, consisting of all the transmission categories, we estimated
181 σ for MSM and HSX using the ratios between σ corresponding to MSM or HSX and the total popu-
182 lation reported from other studies (see footnote in Table S4). Similarly, for obtaining σ for HSX men
183 and women, we employed the corresponding ratios of maximum σ from the CASCADE study. When
184 information necessary to estimate σ was unavailable, we used the highest σ available from the most
185 relatable dataset, as with the UK and the CASCADE study. To estimate the SD of $R_{T/F}$, we employed the

186 error propagation equation²⁴ and derived $\sigma(R_{T/F}) = \frac{\sqrt{\sigma_{infected}^2(\mu_{healthy} - T_{AIDS})^2 + \sigma_{healthy}^2(\mu_{infected} - T_{AIDS})^2}}{(\mu_{healthy} - T_{AIDS})^2}$, where
 187 μ, σ are given in Tables 1 and S6. For $\sigma(R_{T/F})$ of all the data combined, we chose σ from EU/EEA,
 188 involving data from 21 countries.

189 **Estimation of CD4 counts at seroconversion.** In the US study¹⁴, a model of CD4 count decline fol-
 190 lowing seroconversion has been proposed, which allowed us to estimate CD4 counts at seroconversion
 191 from measurements at diagnosis. According to the model, the CD4 count T in an untreated individual
 192 at time t from seroconversion follows $\sqrt{T} = a_0 + b_1 \times t + e_{1t}$, where a_0 and b_1 are constants and e_{1t} is an
 193 error term. At seroconversion, the CD4 count, T_0 , was obtained by setting $t = 0$, so that $\sqrt{T_0} = a_0 + e_{10}$.
 194 Assuming that $e_{1t} = e_{10}$, it followed that $T_0 = (\sqrt{T} - b_1 \times t)^2$. The values of b_1 for different age groups
 195 and transmission categories were available²⁵. Also, the median delays (and IQR) in diagnosis follow-
 196 ing seroconversion, t_d , have been estimated¹⁴, using which we calculated the corresponding mean
 197 and SD. For MSM and HSX, we took the mid-value of the means of t_d in 2006 and 2015 and chose
 198 the largest SD, and obtained $t_d = 4.05 \pm 6.67$ and $t_d = 5.40 \pm 9.04$ years, respectively, for the duration
 199 2006-15. If T_d is the CD4 count at diagnosis, then $T_0 = (\sqrt{T_d} - b_1 \times t_d)^2$. We applied the analysis to
 200 data from the most populated age group (13-29 years) and used the mid-value, 21 years, for which b_1
 201 was $-0.93, -0.77$, and -0.80 year^{-1} for MSM, HSX men, and HSX women, respectively. Furthermore,
 202 we assumed the fractions of females to be the same among all non-MSM groups, in order to obtain
 203 the population sizes of MSM and HSX in this age group (footnote in Table S3). Correspondingly, we
 204 obtained $b_1 = -0.79 \text{ year}^{-1}$ for HSX. To obtain uncertainties in the estimates of T_0 , we repeated the
 205 above analysis with T_d and t_d set at values $\pm \sigma$ away from their respective means, but ensuring that
 206 their lowerbounds ≥ 0 and omitting terms that are second order in σ . Half the difference between the
 207 resulting maximum and minimum values of T_0 yielded the σ corresponding to seroconversion.

208 **Statistical analysis.** To examine whether the mean CD4 counts (or mean $R_{T/F}$) were significantly
 209 higher (or lower) in MSM than HSX, we employed the one-tailed t-test with unequal variance with
 210 the test statistic $t = \frac{\mu_{HSX} - \mu_{MSM}}{\sqrt{\sigma_{HSX}^2/n_{HSX} - \sigma_{MSM}^2/n_{MSM}}}$ and degrees of freedom $d = \frac{\left[\frac{1}{n_{HSX}} + \frac{(\sigma_{MSM}/\sigma_{HSX})^2}{n_{MSM}} \right]^2}{\left[\frac{1}{n_{HSX}^2(n_{HSX}-1)} + \frac{(\sigma_{MSM}/\sigma_{HSX})^4}{n_{MSM}^2(n_{MSM}-1)} \right]}$, where
 211 n_{HSX} and n_{MSM} were the two sample sizes, respectively²⁶. The tests were performed using the R
 212 package²⁷, which yielded corresponding P values.

213 **Data of HIV-1 subtype prevalence.** To assess the extent of mixing between MSM and HSX, we

214 collated data of the prevalence of HIV-1 subtypes in the two groups across relevant geographical
215 regions and calendar years. The data are summarized in Fig. 2 and Tables S1-S2 and discussed in
216 Text S1.

217 **Data of association with transmission clusters.** Finally, we considered the extent of association of
218 MSM and HSX with transmission clusters as an indicator of the time of onward transmission post-
219 infection. The corresponding data we collated along with data of the compositions of the largest
220 transmission clusters in different settings are in Fig. 3 and Table S7 and are discussed in Text S2.

221 **Acknowledgments**

222 We thank Pranesh Padmanabhan, Rajat Desikan, and Pradeep Nagaraja for comments. This work was
223 supported by the DBT/Wellcome Trust India Alliance Senior Fellowship IA/S/14/1/501307 (NMD).

224 **References**

- 225 1 Carlson, J. M. *et al.* Selection bias at the heterosexual HIV-1 transmission bottleneck. *Science* **345**
226 (2014). URL <https://science.sciencemag.org/content/345/6193/1254031>.
- 227 2 Joseph, S. B., Swanstrom, B., Kashuba, A. D. M. & Cohen, M. S. Bottlenecks in HIV-1 transmission:
228 insights from the study of founder viruses. *Nat. Rev. Microbiol.* **13**, 414–425 (2015). URL <https://doi.org/10.1038/nrmicro3471>.
- 229 3 Tully, D. C. *et al.* Differences in the selection bottleneck between modes of sexual transmission
230 influence the genetic composition of the HIV-1 founder virus. *PLOS Pathogens* **12**, 1–29 (2016).
231 URL <https://doi.org/10.1371/journal.ppat.1005619>.
- 232 4 Claiborne, D. T. *et al.* Replicative fitness of transmitted HIV-1 drives acute immune activation,
233 proviral load in memory CD4+ T cells, and disease progression. *Proc. Natl. Acad. Sci. U.S.A.* **112**,
234 E1480–1489 (2015). URL <https://www.pnas.org/content/112/12/E1480>.
- 235 5 Carlson, J. M. *et al.* Impact of pre-adapted HIV transmission. *Nat. Med.* **22(6)**, 606–13 (2016).
236 URL <https://doi.org/10.1038/nm.4100>.
- 237 6 Selhorst, P. *et al.* Replication capacity of viruses from acute infection drives HIV-1 disease
238 progression. *Journal of Virology* **91** (2017). URL [https://jvi.asm.org/content/91/8/
239 e01806-16](https://jvi.asm.org/content/91/8/e01806-16).
- 240 7 Patel, P. *et al.* Estimating per-act HIV transmission risk: a systematic review. *AIDS* **28(10)**, 1509–
241 19 (2014).
242

- 243 8 Owen, B. N. *et al.* Prevalence and frequency of heterosexual anal intercourse among young people:
244 a systematic review and meta-analysis. *AIDS Behav.* **19(7)**, 1338–60 (2015). URL <https://doi.org/10.1007/s10461-015-0997-y>.
245
- 246 9 Deeks, S. G., Overbaugh, J., Phillips, A. & Buchbinder, S. HIV infection. *Nat. Rev. Dis. Primers* **1**,
247 15035 (2015).
- 248 10 European Centre for Disease Prevention and Control. *HIV/AIDS surveillance in Europe 2019*
249 (HIV/AIDS surveillance in Europe 2019 - 2018 data). URL <https://www.ecdc.europa.eu/en/publications-data/hivaids-surveillance-europe-2019-2018-data>. [Online;
250 accessed 01-August-2020].
251
- 252 11 Frentz, D. *et al.* Patterns of transmitted HIV drug resistance in Europe vary by risk group. *PLoS*
253 *One* **9(4)**, e94495 (2014). URL <https://doi.org/10.1371/journal.pone.0094495>.
- 254 12 Gupta, S. B., Gilbert, R. L., Brady, A. R., Livingstone, S. J. & Evans, B. G. CD4 cell counts in adults
255 with newly diagnosed HIV infection: results of surveillance in England and Wales, 1990-1998.
256 *AIDS* **14(7)**, 853–861 (2000).
- 257 13 Tang, H. *et al.* Baseline CD4 cell counts of newly diagnosed HIV cases in China: 2006–2012. *PLoS*
258 *ONE* **9**, e96098 (2014). URL <https://doi.org/10.1371/journal.pone.0096098>.
- 259 14 Robertson, M. M., Braunstein, S. L., Hoover, D. R., Li, S. & Nash, D. Estimates of the time from
260 seroconversion to ART initiation among people newly diagnosed with HIV from 2006 to 2015,
261 New York City. *Clin. Infect. Dis.* **ciz1178** (2019).
- 262 15 CASCADE Collaboration. Differences in CD4 cell counts at seroconversion and decline among
263 5739 HIV-1-infected individuals with well-estimated dates of seroconversion. *J. Acquir. Im-*
264 *mune Defic. Syndr.* **34**, 76–83 (2003). URL [https://journals.lww.com/jaids/Fulltext/](https://journals.lww.com/jaids/Fulltext/2003/09010/Differences_in_CD4_Cell_Counts_at_Seroconversion.12.aspx)
265 [2003/09010/Differences_in_CD4_Cell_Counts_at_Seroconversion.12.aspx](https://journals.lww.com/jaids/Fulltext/2003/09010/Differences_in_CD4_Cell_Counts_at_Seroconversion.12.aspx).
- 266 16 Lodi, S. *et al.* Time from human immunodeficiency virus seroconversion to reaching CD4+ cell
267 count thresholds < 200 , < 350 , and < 500 cells/mm³ : assessment of need following changes in
268 treatment guidelines. *Clin. Infect. Dis.* **53(8)**, 817–825 (2011). URL [https://doi.org/10.](https://doi.org/10.1093/cid/cir494)
269 [1093/cid/cir494](https://doi.org/10.1093/cid/cir494).
- 270 17 Beyrer, C. *et al.* Global epidemiology of HIV infection in men who have sex with men.
271 *Lancet (London, England)* **380(9839)**, 367–377 (2012). URL <https://doi.org/10.1016/>

- 272 S0140-6736 (12) 60821-6.
- 273 18 Ariën, K. K., Vanham, G. & Arts, E. J. Is HIV-1 evolving to a less virulent form in humans? *Nat.*
274 *Rev. Microbiol.* **5**, 141–151 (2007). URL <https://doi.org/10.1038/nrmicro1594>.
- 275 19 Shet, A., Nagaraja, P. & Dixit, N. M. Viral decay dynamics and mathematical modeling of treatment
276 response: evidence of lower in vivo fitness of HIV-1 subtype C. *J. Acquir. Immune Defic. Syndr.*
277 **73**, 245–251 (2016). URL [https://journals.lww.com/jaids/Fulltext/2016/11010/
278 Viral_Decay_Dynamics_and_Mathematical_Modeling_of.1.aspx](https://journals.lww.com/jaids/Fulltext/2016/11010/Viral_Decay_Dynamics_and_Mathematical_Modeling_of.1.aspx).
- 279 20 Villabona-Arenas, C. J. *et al.* Number of HIV-1 founder variants is determined by the recency of the
280 source partner infection. *Science* **369**, 103–108 (2020). URL [https://science.sciencemag.
281 org/content/369/6499/103](https://science.sciencemag.org/content/369/6499/103).
- 282 21 Fraser, C. *et al.* Virulence and pathogenesis of HIV-1 infection: An evolutionary perspective. *Science*
283 **343** (2014). URL <https://science.sciencemag.org/content/343/6177/1243727>.
- 284 22 Wan, X., Wang, W., Liu, J. & Tong, T. Estimating the sample mean and standard deviation from
285 the sample size, median, range and/or interquartile range. *BMC Med. Res. Methodol.* **14** (2014).
286 URL <https://doi.org/10.1186/1471-2288-14-135>.
- 287 23 Higgins J. P. T., Thomas J., Chandler J., Cumpston M., Li T., Page M. J. & Welch V. A.
288 (ed.) *Cochrane Handbook for Systematic Reviews of Interventions* (Cochrane, 2019). URL [www.
289 training.cochrane.org/handbook](http://www.training.cochrane.org/handbook). [version 6.0 (updated July 2019)].
- 290 24 Bevington, P. R. and Robinson, D. K. *Data Reduction and Error Analysis for the Physical Sciences*
291 (McGraw Hill, 2003), 3 edn.
- 292 25 Song, R., Hall, H. I., Green, T. A., Szwarcwald, C. L. & Pantazis, N. Using CD4 data to estimate
293 HIV incidence, prevalence, and percent of undiagnosed infections in the United States. *J. Acquir.*
294 *Immune Defic. Syndr.* **74**, 3–9 (2017).
- 295 26 Ruxton, Graeme D. The unequal variance t-test is an underused alternative to Student's t-test
296 and the Mann–Whitney U test. *Behav. Ecol.* **17**, 688–690 (2006). URL [https://doi.org/10.
297 1093/beheco/ark016](https://doi.org/10.1093/beheco/ark016).
- 298 27 R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical
299 Computing, Vienna, Austria (2013). URL <http://www.R-project.org/>.

SUPPLEMENTARY INFORMATION

Transmitted HIV-1 is more pathogenic in heterosexual individuals than homosexual men

Ananthu James^a and Narendra M. Dixit^{a,b,*}

^aDepartment of Chemical Engineering, Indian Institute of Science, Bengaluru, India

^bCentre for Biosystems Science and Engineering, Indian Institute of Science, Bengaluru, India

*Corresponding author; E-mail: narendra@iisc.ac.in

Contents:

Text notes: 2

Figures: 1

Tables: 7

References: 41

1 **Text S1: Distinct subtype prevalences indicate minimal mixing between MSM and HSX**

2 In many geographical locations, mixing between MSM and HSX appears minimal. This is evident
3 from the different prevalences of HIV-1 subtypes in the two groups. MSM in western nations are
4 dominated by HIV-1 subtype B, whereas HSX comprise a mixture of subtypes¹, with subtypes B and
5 C being the predominant ones². For instance, in the United Kingdom, from 2002-2010, MSM had
6 nearly 90% subtype B infections, whereas HSX had a little over 10% subtype B. Mixing between the
7 two groups would have led to a more similar distribution of subtypes in the two. The two groups thus
8 appear to have sustained their respective infections over the years in near complete isolation. The
9 difference in subtype prevalences holds also in Canada, Spain, France, and other nations (Fig. 2(a);
10 Table S1). In China, the dominant subtype is CRF01-AE, which is present in MSM with a frequency of
11 >50% but in HSX at <40% (Fig. 2(b); Table S2)³, perhaps indicative of more mixing than in Europe.
12 In Korea, the extent of mixing could not be assessed using subtypes because over 80% of all infections
13 were subtype B⁴. In USA, though subtype B dominates both MSM and HSX^{5,6}, mixing between
14 the groups has been argued not to be common⁷. In the Nordic states, some mixing between MSM
15 and HSX is evident⁸. Overall, little mixing between MSM and HSX is evident in most geographical
16 settings, suggesting that the different selection biases between the groups may have been sustained
17 over the course of the epidemic.

19 **Text S2: Clustering and transmission patterns**

20 MSM are known to engage in different sexual contact patterns compared to HSX. They tend to have
21 more partners than HSX^{9,10}. They are also far more likely to belong to transmission clusters compared
22 to HSX¹. A transmission cluster comprises individuals carrying viral genomes that cluster together in
23 a phylogenetic tree¹¹, suggesting that the viral sequences isolated from the individuals are closely
24 related. In Japan and China, an infected MSM had a nearly 40% chance of being part of a cluster,
25 whereas an infected HSX had <10% chance¹². In France, the corresponding numbers were ~35%
26 and ~4%, respectively¹³. This trend was true for all the countries with data available except the
27 Netherlands (Fig. 3(a); Table S7). MSM also formed larger clusters than HSX. The largest clusters
28 reported in Belgium and Spain comprised nearly 100 individuals each, with the Belgian cluster con-
29 taining ~70 MSM and the Spanish cluster exclusively MSM (Fig. 3(b); Table S7)^{14,15}. Together, these

30 data suggest greater similarity in the viral strains in MSM than HSX. One way in which this greater
31 similarity could arise is by onward transmission occurring sooner after infection in MSM than HSX,
32 allowing lesser individual host-specific adaptation before transmission.

33 **Supplementary Figures**

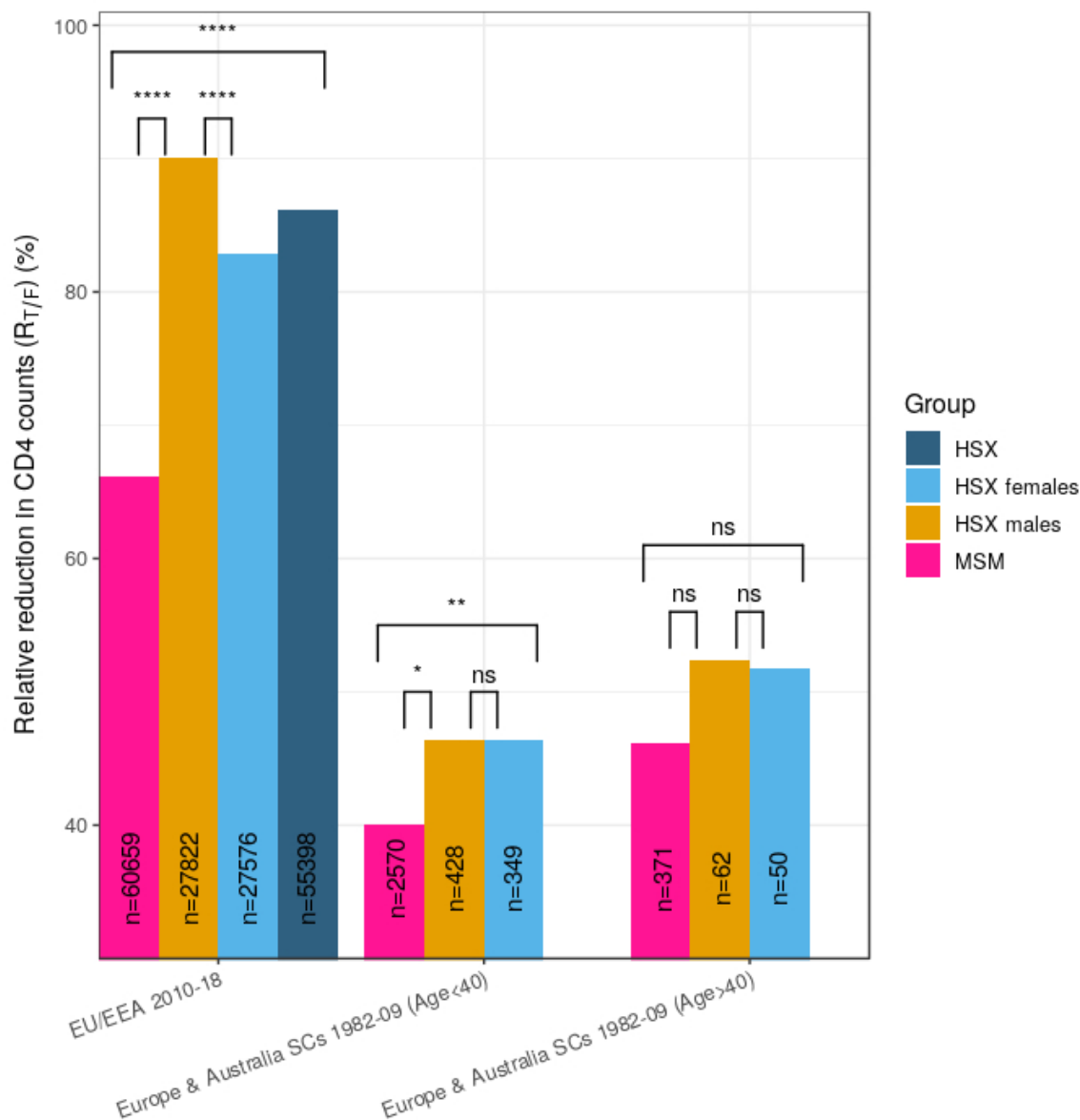


Figure S1. Effect of gender on the relative reduction in early CD4 counts. $R_{T/F}$ from EU/EEA during 2010-18¹⁶ and the CASCADE study¹⁷ indicate that HSX men have higher $R_{T/F}$ than MSM, with this difference achieving significance with large sample sizes, ruling out gender as a cause of the lower CD4 count reduction in MSM than HSX. The sample sizes (n) are indicated along with the P values. (****, ***, ** and * indicate $P < 10^{-4}$, $P < 10^{-3}$, $P < 10^{-2}$ and $P < 0.05$, respectively, while ns (not significant) implies $P > 0.05$.) Additional details are in Table 1 (main text).

34 **Supplementary Tables**

Table S1. Prevalence of subtype B in Europe. Data from different regions in Europe show substantially higher subtype B percentage prevalence in MSM than HSX ($P < 0.001$ in each study, unless specified). The sample sizes (n) are in parantheses.

Region [Reference]	Survey year(s)	Prevalence of subtype B (%)	
		MSM (n)	HSX (n)
France ¹⁸	2001	80.5 (113)	53.2 (186)
	2006 – 07	90.6 (223)	27.2 (202)
	2010 – 11	83.3 (311)	26.2 (263)
Italy ¹⁹	2000 – 14	81.5 (1,824)	54.9 (1,573)
UK ²⁰	2002 – 10	88.4 [‡] (21,321)	13.3 [‡] (38,948)
Madrid, Spain ²¹	2000 – 11	88.2 (340)	41.5 (369)
Comunidad Valenciana (CV), Spain ¹⁵	2004 – 14	92.0 [‡] (637)	60.8 [‡] (204)
Europe ²²	2002 – 07	90.4 (2,084)	33.5 (1,501)
Canada* & Europe ²³	1998 – 10	95.5 [‡] (5,570)	38.7 [‡] (2,052)
Japan ¹²	2004 – 11	96.9 [‡] (261)	39.8 [‡] (103)

[‡] P value not specified.

*In this study, 15.9% patients were from Canada.

Table S2. Prevalence of subtypes in China. A recent review³ of 130 published articles, together involving of 10,516 MSM and 6,759 HSX individuals, has examined the prevalence of different subtypes in China, which is reproduced below. P values indicate significant differences in the prevalences of 3 subtypes.

Subtype	Prevalence (%) (95% CI)		
	MSM	HSX	P value
CRF01_AE	51.28 (46.15 – 56.40)	38.78 (33.08 – 44.63)	< 0.01
CRF07_BC	19.98 (16.17 – 24.07)	14.88 (10.96 – 19.23)	0.083
CRF08_BC	0.00 (0.00 – 0.00)	9.81 (6.44 – 13.70)	< 0.01
B\B'	17.74 (12.78 – 23.26)	15.41 (11.15 – 20.16)	0.508
C	0.00 (0.00 – 0.00)	1.89 (0.81 – 3.28)	< 0.01
Others	1.56 (0.98 – 2.24)	2.39 (1.33 – 3.66)	0.2109

Table S3. Early median CD4 cell counts in infected adults from several large population studies*. The sources of the studies, the periods of study, measurement times, and other details are mentioned. The USA and European studies report IQRs, whereas the CASCADE study provides 95% CIs. NA - not available.

Region	Time period	Measurement time	Risk group	CD4 counts (cells/ μ L)	Sample size (n)	Other details
USA ^{‡24}	2006 – 15	Within 6 m.o. diagnosis	MSM	387 (227 – 555)	17,779	
			HSX	308 (125 – 508)	8,310	
Europe ^{†22}	2002 – 07	Within 6 m.o. diagnosis	MSM	435 (259 – 585)	2,084	$P < 10^{-3}$
			HSX	280 (110 – 458)	1,501	
England & Wales ²⁵	1990 – 98	Within 6 m.o. diagnosis	MSM	331 (NA)**	6,213	$P < 10^{-3}$
			HSX	230 (NA)**	2,637	
Europe & Australia ¹⁷ (CASCADE) SCs ^{††}	1979 – 00	Seroconversion	MSM (age < 40)	621 (609 – 634)	2,941	n=5739. 54.7% (MSM & hemophilic men) only men. 20.4% women. 12.6% aged > 40.
			MSM (age > 40)	578 (556 – 601)		
			HSX men (age < 40)	576 (545 – 607)	889	
			HSX men (age > 40)	534 (501 – 569)		
			HSX women (age < 40)	623 (599 – 647)		
			HSX women (age > 40)	580 (548 – 612)		

*We did not include studies comparing CD4 counts in MSM and the general population, the latter consisting of not only MSM and HSX, but also considerable fractions of injection drug users and hemophiliacs²⁶. The latter two groups may have significantly different early CD4 counts than MSM and HSX¹⁷, precluding a direct comparison between MSM and HSX. When many articles were available for overlapping durations, corresponding to the same geographic location (as in the case of New York City), we chose data from the study having a higher sample size²⁴. Interestingly, we found no study with $n \gtrsim 1000$ where mean/median early cell counts among HSX were significantly higher than those among MSM⁴.

[‡]Corresponding to this population, the proportion of females was 22.2% and the age group 13-29 had 35.6% of $n = 28,162$ individuals.

[†]Median viral load (with IQR) in log copies/mL is 4.88 (4.3 – 5.4) for MSM and 4.79 (4.2 – 5.3) for HSX ($P < 0.001$)²².

**We chose the SDs from the Europe (2002-07) study (see Methods).

^{††}For the whole sets of MSM and HSX, including all the age groups, we used the SDs corresponding to age < 40 years. For HSX, we used the largest SD among HSX men and women.

Table S4. Early median CD4 cell counts in infected MSM and HSX from China²⁷. The sample sizes (n) were directly available, while the cell counts were estimated using WebPlotDigitizer²⁸. The last row represents estimates (see Methods) for the entire period 2006-12.

Year	MSM		HSX	
	Cell counts (<i>cells/μL</i>)	Sample size (n)	Cell counts (<i>cells/μL</i>)	Sample size (n)
2006	315	244	230	3,305
2007	345	672	249	7,251
2008	362	2,025	266	12,378
2009	364	3,624	265	18,208
2010	368	5,648	267	24,425
2011	370	9,316	271	35,211
2012	371	13,748	281	42,653
Average	368 ± 222 (SD [†])	35,277	270 ± 260 (SD [†])	143,431

[†]Corresponding to the total population that contributed to the CD4 count data in China, consisting of MSM, HSX and other transmission categories, during 2006 – 12, the IQR was 130 – 454 *cells/μL*. Using this, we obtained $SD = 240$ *cells/μL* (see Methods). To calculate the SDs corresponding to MSM and HSX, we recall that in the USA²⁴ the SDs corresponding to the CD4 counts from MSM, HSX and the total population were 243, 284 and 262 *cells/μL*, respectively. Assuming that the ratios between SDs remain the same in China as well, we evaluated the SD for MSM to be 222 and HSX to be 260 *cells/μL* in China.

Table S5. Early median CD4 cell counts in infected adults from EU/EEA¹⁶. The cell counts were estimated using WebPlotDigitizer²⁸. The median ages, whenever available, are also provided. The last row provides the mean cell counts (with SDs) and total numbers of MSM, HSX men and women, respectively, estimated as in Methods.

Year	No. of diagnoses*		Fraction with CD4 counts ^{††}	Median age		Median CD4 count (<i>cells/μL</i>)			
	MSM	HSX		MSM	HSX	MSM	HSX men	HSX women	
2010	10,348	11,693	0.591	-	-	420	261	312	
2011	10,411	11,011	0.561	-	-	422	260	317	
2012	11,238	10,782	0.553	-	-	436	270	319	
2013	11,553	10,187	0.607	-	-	443	271	339	
2014	11,821	10,049	0.614	-	-	450	281	346	
2015	11,484	9,267	0.747	-	-	454	284	345	
2016	10,508	9,003	0.673	34	39	440	274	351	
2017	9,723	8,518	0.717	34	39	434	268	349	
2018	8,049	7,267	0.695	36	41	423	262	340	
Total	$n_{MSM}^{**} = 60,659, n_{HSXM}^{**} = 27,822, n_{HSXW}^{**} = 27,576$						$(437 \pm 397)^{\ddagger}$	$(270 \pm 464)^{\ddagger}$	$(335 \pm 325)^{\ddagger}$

*The number of diagnoses specific to each country during 2009-18 was provided for MSM and HSX in the latest annual report¹⁶. Using these, we estimated the numbers above.

††These are taken from the annual surveillance reports²⁹. Since this number was unavailable for 2009, we omitted the data for the calculation of the mean CD4 counts for the combined data.

**The numbers of HSX men and women diagnosed from the 21 EU/EEA countries in 2018 were available in the 2019 annual report¹⁶. Using this, we calculated the fractions of HSX men and women and assumed that they remained the same during 2010-18. Also, $n_{HSX} = n_{HSXmen} + n_{HSXwomen} = 55,398$. (There were a small number of transgenders among HSX in column 3, but the CD4 counts were reported only for HSX men and women.)

‡In 2018, for 16,694 adults across risk groups, the median cell counts with 95% CIs were 365 (359 – 372) *cells/μL*. Using this, we computed their SD to be ~ 429 *cells/μL* and evaluated the SDs for MSM and HSX. To obtain the SDs for HSX men and women, we used the SD ratio for the age group < 40 years corresponding to the CASCADE study (see Methods).

Table S6. CD4 T cell counts in healthy adults. Mean CD4 counts in healthy adults from different population groups which define baseline counts for estimating the relative reduction in early cell count following HIV-1 infection. Sample sizes are in brackets. SD is standard deviation.

Region	Category (sample size, <i>n</i>)	Cell counts (<i>cells/μL</i>)	Other details
China ^{30,31} (Hong Kong ³⁰ & Shanghai ³¹)	Men (78 ³⁰ & 377 ³¹)	725	Combined mean; SD = 258 ³⁰ & 256 ³¹ .
	Women (130 ³⁰ & 237 ³¹)	724	Combined mean; SD = 254 ³⁰ & 255 ³¹ .
	HSX (822)	725	SD = 255 for Shanghai ³¹
Sub-Saharan Africa (Tanzania) ³²	Men* (42)	666	SD = 247
	Women* (60)	802	SD = 250
	HSX (102)	746	SD = 257 ^{†‡}
USA ³³	Men** (33)	921	SD = 188
	Women (67)	1041	SD = 340
	HSX** (100)	1001	SD = 305
UK ³⁴	HSX men [†] (50)	840	SD = 285
	HSX women (50)	1050	SD = 377
	HSX [‡] (100)	945	We used the SD for women.
	MSM ^{†,‡} (100)	800	SD = 324
EU/EEA (& Europe) & Australia (Italy) ³⁵	General population or HSX ^{††} (965)	941	SDs were not reported. We used the SDs from the UK study (see Methods).
	Men ^{††} (532)	902	
	Women (436)	989	

* $P = 8 \times 10^{-3}$ for the comparison between men and women (reported in the original study³²).

** $P = 0.04$ for the comparison between healthy men and HSX from USA.

[†] $P = 0.22$ for the comparison between MSM and HSX men from UK.

[‡] $P = 6.6 \times 10^{-4}$ for the comparison between healthy MSM and HSX from UK, using SD (= 305) for HSX from USA. Using SD (=377) corresponding to HSX women from UK instead yielded $P = 2 \times 10^{-3}$.

^{††} $P = 0.018$ using the SDs from UK.

^{‡‡}For our calculation involving the EU/EEA 2010-18 population (footnote in Table 1), we used the SD (=377) for the HSX from Europe/UK.

Table S7. Association of MSM and HSX with transmission clusters. The percentages of individuals found to be associated with transmission clusters in MSM and HSX in several studies are collated. In the second column, the numbers in parentheses indicate sample sizes examined. Where available, P values and the largest cluster sizes are indicated in other details.

Region (Study period)	Group (n)	Associated with clusters (%)	Other details
Japan ¹² (2004 – 11)	MSM (261)	38.3	Clusters are also linked to China.
	HSX (103)	8.7	
Austria ³⁶ (2008 – 14)	MSM (122)	57.4	$P < 0.001$
	HSX (115)	32.2	
Belgium ¹⁴ (2013 – 15)	MSM (698)	57.7	$P < 0.001$. Cluster size ≥ 3 . The largest cluster is sized 94, consisting of 68 MSM and 15 HSX.
	HSX (508)	22.4	
Cyprus ³⁷ (1986 – 2012)	MSM (164)	36.6	Cluster size ≥ 3 . The two largest clusters are sized 13, having 11 – 12 MSM.
	HSX (129)	2.3	
France ¹³ (1999 – 2014)	MSM (636)	32.1	$P < 0.01$. Cluster size ≥ 3 . The largest cluster, sized 41, has 39 MSM and a HSX.
	HSX (267)	3.8	
Germany ³⁸ (1999 – 2016)	MSM (1,448)	27.3	
	HSX (622)	18.2	
Netherlands ³⁹ (1981 – 2011)	MSM (4,288)	51.2	Cluster size ≥ 10 .
	HSX (849)	55.8	
Spain ¹⁵ (2004 – 14)	MSM (637)	–	8 of 12 clusters sized ≥ 10 consist only of MSM. The largest cluster has a size of 111, and is composed only of MSM.
	HSX (204)	–	
USA ⁴⁰ (2001 – 15)	MSM (1,597)	37.4	The largest cluster has a size of 39, and has only MSM.
	HSX (953)	28.3	
Canada and 9 European countries ⁴¹ (1981 – 2011)	MSM (4,980)	51.4	The 9 countries are Italy, Greece, Netherlands, Germany, France, UK, Norway, Austria, and Spain. $P < 0.001$.
	Bisexual men (2,087)	27.7	

35 References

- 36 1 Beyrer, C., Baral, S. D., van Griensven, F., Goodreau, S. M., Chariyalertsak, S., Wirtz, A. L.,
37 and Brookmeyer, R. Global epidemiology of HIV infection in men who have sex with men.
38 *Lancet (London, England)* **380(9839)**, 367–377 (2012). URL [https://doi.org/10.1016/](https://doi.org/10.1016/S0140-6736(12)60821-6)
39 [S0140-6736\(12\)60821-6](https://doi.org/10.1016/S0140-6736(12)60821-6).
- 40 2 Abecasis, A. B. et al. HIV-1 subtype distribution and its demographic determinants in newly diag-
41 nosed patients in Europe suggest highly compartmentalized epidemics. *Retrovirology* **10** (2013).
42 URL <https://doi.org/10.1186/1742-4690-10-7>.
- 43 3 Yuan, R., Cheng, H., Chen, L. S., Zhang, X., and Wang, B. Prevalence of different HIV-1 subtypes in
44 sexual transmission in China: a systematic review and meta-analysis. *Epidemiol. Infect.* **144(10)**,
45 2144–2153 (2016).
- 46 4 Kim, G. J. et al. Estimating the origin and evolution characteristics for Korean HIV type 1 subtype
47 B using Bayesian phylogenetic analysis. *AIDS Res. Hum. Retroviruses* **28(8)**, 880–884 (2012). URL
48 <https://doi.org/10.1089/aid.2011.0267>.
- 49 5 Junqueira, D. M. and Almeida, S. E. d. M. HIV-1 subtype B: traces of a pandemic. *Virology* **495**,
50 173–184 (2016).
- 51 6 Dennis, A. M. et al. Rising prevalence of non-B HIV-1 subtypes in North Carolina and evidence
52 for local onward transmission. *Virus Evol.* **3(1)**, vex013 (2017). URL [https://doi.org/10.](https://doi.org/10.1093/ve/vex013)
53 [1093/ve/vex013](https://doi.org/10.1093/ve/vex013).
- 54 7 Satcher, A. J., Durant, T., Hu, X., and Dean, H. D. AIDS cases among women who reported sex
55 with a bisexual man, 2000-2004—United States. *Women & Health* **46**, 23–40 (2007).
- 56 8 Esbjörnsson, J. et al. HIV-1 transmission between MSM and heterosexuals, and increasing propor-
57 tions of circulating recombinant forms in the Nordic Countries. *Virus Evolution* **2** (2016). URL
58 <https://doi.org/10.1093/ve/vew010>.
- 59 9 Glick, S. N. et al. A comparison of sexual behavior patterns among men who have sex with men
60 and heterosexual men and women. *J. Acquir. Immune Defic. Syndr.* **60(1)**, 83–90 (2012). URL
61 <https://doi.org/10.1097/QAI.0b013e318247925e>.
- 62 10 Kenyon, C. R., Wolfs, K., Osbak, K., van Lankveld, J., and Van Hal, G. Implicit attitudes to sexual
63 partner concurrency vary by sexual orientation but not by gender—A cross sectional study of

- 64 Belgian students. *PLoS One* **13(5)**, e0196821 (2018).
- 65 11 Hassan, A. S., Pybus, O. G., Sanders, E. J., Albert, J., and Esbjörnsson, J. Defining HIV-1 transmis-
66 sion clusters based on sequence data. *AIDS* **31(9)**, 1211–1222 (2017).
- 67 12 Kondo, M. *et al.* Emergence in Japan of an HIV-1 variant associated with transmission among men
68 who have sex with men (MSM) in China: first indication of the international dissemination of
69 the Chinese MSM lineage. *Journal of Virology* **87**, 5351–5361 (2013). URL [https://jvi.asm.
70 org/content/87/10/5351](https://jvi.asm.org/content/87/10/5351).
- 71 13 Chaillon, A. *et al.* Spatiotemporal dynamics of HIV-1 transmission in France (1999-2014) and
72 impact of targeted prevention strategies. *Retrovirology* **14(1)**, 15 (2017). URL [https://doi.
73 org/10.1186/s12977-017-0339-4](https://doi.org/10.1186/s12977-017-0339-4).
- 74 14 Verhofstede, C. *et al.* Phylogenetic analysis of the Belgian HIV-1 epidemic reveals that lo-
75 cal transmission is almost exclusively driven by men having sex with men despite presence
76 of large African migrant communities. *Infect. Genet. Evol.* **61**, 36–44 (2018). URL [https:
77 //doi.org/10.1016/j.meegid.2018.03.002](https://doi.org/10.1016/j.meegid.2018.03.002).
- 78 15 Patino-Galindo, J. A. *et al.* The molecular epidemiology of HIV-1 in the Comunidad Valenciana
79 (Spain): analysis of transmission clusters. *Sci. Rep.* **7**, 11584 (2017). URL [https://doi.org/
80 10.1038/s41598-017-10286-1](https://doi.org/10.1038/s41598-017-10286-1).
- 81 16 European Centre for Disease Prevention and Control. *HIV/AIDS surveillance in Europe 2019*
82 (HIV/AIDS surveillance in Europe 2019 - 2018 data). URL [https://www.ecdc.europa.eu/
83 en/publications-data/hivaids-surveillance-europe-2019-2018-data](https://www.ecdc.europa.eu/en/publications-data/hivaids-surveillance-europe-2019-2018-data). [Online;
84 accessed 01-August-2020].
- 85 17 CASCADE Collaboration. Differences in CD4 cell counts at seroconversion and decline among
86 5739 HIV-1-infected individuals with well-estimated dates of seroconversion. *J. Acquir. Im-
87 mune Defic. Syndr.* **34**, 76–83 (2003). URL [https://journals.lww.com/jaids/Fulltext/
88 2003/09010/Differences_in_CD4_Cell_Counts_at_Seroconversion.12.aspx](https://journals.lww.com/jaids/Fulltext/2003/09010/Differences_in_CD4_Cell_Counts_at_Seroconversion.12.aspx).
- 89 18 Descamps, D. *et al.* National sentinel surveillance of transmitted drug resistance in antiretroviral-
90 naive chronically HIV-infected patients in France over a decade: 2001–2011. *J. Antimicrob.
91 Chemother.* **68(11)**, 2626–2631 (2013). URL <https://doi.org/10.1093/jac/dkt238>.
- 92 19 Fabeni, L. *et al.* Dynamics and phylogenetic relationships of HIV-1 transmitted drug resistance ac-

- 93 cording to subtype in Italy over the years 2000–14. *J. Antimicrob. Chemother.* **72(10)**, 2837–2845
94 (2017). URL <https://doi.org/10.1093/jac/dkx231>.
- 95 20 The UK Collaborative Group on HIV Drug Resistance. The increasing genetic diver-
96 sity of HIV-1 in the UK, 2002–2010. *AIDS* **28(5)**, 773–780 (2014). URL [https://journals.lww.com/aidsonline/Fulltext/2014/03130/The_increasing_](https://journals.lww.com/aidsonline/Fulltext/2014/03130/The_increasing_genetic_diversity_of_HIV_1_in_the.15.aspx)
97 [genetic_diversity_of_HIV_1_in_the.15.aspx](https://journals.lww.com/aidsonline/Fulltext/2014/03130/The_increasing_genetic_diversity_of_HIV_1_in_the.15.aspx).
- 98
99 21 Yebra, G. et al. Different trends of transmitted HIV-1 drug resistance in Madrid, Spain, among risk
100 groups in the last decade. *Arch. Virol.* **159(5)**, 1079–87 (2014). URL [https://doi.org/10.](https://doi.org/10.1007/s00705-013-1933-y)
101 [1007/s00705-013-1933-y](https://doi.org/10.1007/s00705-013-1933-y).
- 102 22 Frentz, D. et al. Patterns of transmitted HIV drug resistance in Europe vary by risk group. *PLoS*
103 *One* **9(4)**, e94495 (2014). URL <https://doi.org/10.1371/journal.pone.0094495>.
- 104 23 Klein, M. B. et al. The effects of HIV-1 subtype and ethnicity on the rate of CD4 cell count decline
105 in patients naive to antiretroviral therapy: a Canadian-European collaborative retrospective co-
106 hort study. *CMAJ OPEN* **2(4)**, E318–E329 (2014). URL [https://doi.org/10.9778/cmajo.](https://doi.org/10.9778/cmajo.20140017)
107 [20140017](https://doi.org/10.9778/cmajo.20140017).
- 108 24 Robertson, M. M., Braunstein, S. L., Hoover, D. R., Li, S., and Nash, D. Estimates of the time from
109 seroconversion to ART initiation among people newly diagnosed with HIV from 2006 to 2015,
110 New York City. *Clin. Infect. Dis.* **ciz1178** (2019).
- 111 25 Gupta, S. B., Gilbert, R. L., Brady, A. R., Livingstone, S. J., and Evans, B. G. CD4 cell counts in
112 adults with newly diagnosed HIV infection: results of surveillance in England and Wales, 1990-
113 1998. *AIDS* **14(7)**, 853–861 (2000).
- 114 26 Pantazis, N. et al. Temporal trends in prognostic markers of HIV-1 virulence and transmissibility:
115 an observational cohort study. *Lancet HIV* **1(3)**, e119–26 (2014). URL [https://doi.org/10.](https://doi.org/10.1016/S2352-3018(14)00002-2)
116 [1016/S2352-3018\(14\)00002-2](https://doi.org/10.1016/S2352-3018(14)00002-2).
- 117 27 Tang, H. et al. Baseline CD4 cell counts of newly diagnosed HIV cases in China: 2006–2012. *PLoS*
118 *ONE* **9**, e96098 (2014). URL <https://doi.org/10.1371/journal.pone.0096098>.
- 119 28 Rohatgi, A. *WebPlotDigitizer*. Pacifica, California, USA (July, 2020). URL [https://automeris.](https://automeris.io/WebPlotDigitizer)
120 [io/WebPlotDigitizer](https://automeris.io/WebPlotDigitizer).
- 121 29 European Centre for Disease Prevention and Control. *Annual HIV/AIDS surveillance re-*

- ports (HIV/AIDS surveillance in Europe). URL <https://www.ecdc.europa.eu/en/all-topics-zhiv-infection-and-aids-surveillance-and-disease-data/annual-hiv-aids-surveillance-reports>. [Online; accessed 01-August-2020].
- 30 Kam, K. M. *et al.* Lymphocyte subpopulation reference ranges for monitoring human immunodeficiency virus-infected Chinese adults. *Clinical and Vaccine Immunology* **3**, 326–330 (1996). URL <https://cvi.asm.org/content/3/3/326>.
- 31 Jiang, W. *et al.* Normal values for CD4 and CD8 lymphocyte subsets in healthy Chinese adults from Shanghai. *Clinical and Vaccine Immunology* **11**, 811–813 (2004). URL <https://cvi.asm.org/content/11/4/811>.
- 32 Ngowi, B. J., Mfinanga, S. G., Bruun, J. N., and Morkve, O. Immunohaematological reference values in human immunodeficiency virus-negative adolescent and adults in rural northern Tanzania. *BMC Infect Dis.* **9** (2009).
- 33 Valiathan, R. *et al.* Reference ranges of lymphocyte subsets in healthy adults and adolescents with special mention of T cell maturation subsets in adults of South Florida. *Immunobiology* **219**(7), 487–496 (2014). URL <https://doi.org/10.1016/j.imbio.2014.02.010>.
- 34 Bofill, M. e. a. Laboratory control values for CD4 and CD8 T lymphocytes. Implications for HIV-1 diagnosis. *Clinical & Experimental Immunology* **88**(2), 243–252 (1992). URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1365-2249.1992.tb03068.x>.
- 35 Santagostino, A. *et al.* An Italian national multicenter study for the definition of a reference ranges for normal values of peripheral blood lymphocyte subsets in healthy adults. *Haematologica* **84**(6), 499–504 (1999). URL <https://pubmed.ncbi.nlm.nih.gov/10366792/>.
- 36 Hoenigl, M. *et al.* Characterization of HIV transmission in South-east Austria. *PLoS ONE* **11**(3), e0151478 (2016). URL <https://doi.org/10.1371/journal.pone.0151478>.
- 37 Pineda-Peña, A. *et al.* HIV-1 infection in Cyprus, the Eastern Mediterranean European frontier: a densely sampled transmission dynamics analysis from 1986 to 2012. *Sci. Rep.* **8** (2018). URL <https://doi.org/10.1038/s41598-017-19080-5>.
- 38 Stecher, M. *et al.* Molecular epidemiology of the HIV epidemic in three German metropolitan regions – Cologne/Bonn, Munich and Hannover, 1999–2016. *Sci. Rep.* **8** (2018). URL <https://doi.org/10.1038/s41598-018-25004-8>.

- 151 39 Bezemer, D. *et al.* Dispersion of the HIV-1 epidemic in men who have sex with men in the Nether-
152 lands: a combined mathematical model and phylogenetic analysis. *PLoS Med.* **12(11)**, e1001898
153 (2015). URL <https://doi.org/10.1371/journal.pmed.1001898>.
- 154 40 Dennis, A. M. *et al.* HIV-1 transmission clustering and phylodynamics highlight the important role
155 of young men who have sex with men. *AIDS Res. Hum. Retroviruses* **34(10)**, 879–88 (2018). URL
156 <https://doi.org/10.1089/aid.2018.0039>.
- 157 41 Paraskevis, D. *et al.* HIV-1 molecular transmission clusters in nine European countries and Canada:
158 association with demographic and clinical factors. *BMC Medicine* **17** (2019). URL [https://doi.](https://doi.org/10.1186/s12916-018-1241-1)
159 [org/10.1186/s12916-018-1241-1](https://doi.org/10.1186/s12916-018-1241-1).