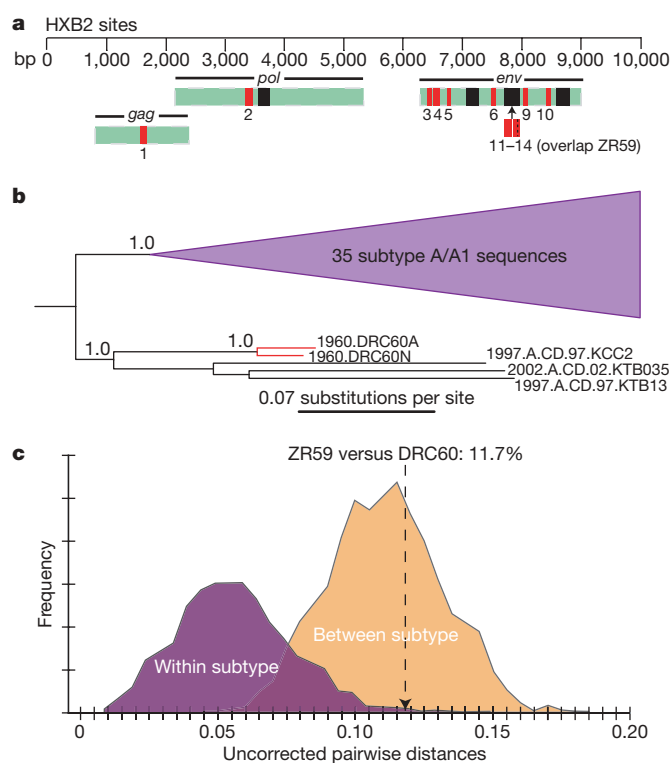


# Direct evidence of extensive diversity of HIV-1 in Kinshasa by 1960

Michael Worobey<sup>1</sup>, Marlea Gemmel<sup>1</sup>, Dirk E. Teuwen<sup>2,3</sup>, Tamara Haselkorn<sup>1</sup>, Kevin Kunstman<sup>4</sup>, Michael Bunce<sup>5</sup>, Jean-Jacques Muyembe<sup>6,7</sup>, Jean-Marie M. Kabongo<sup>6</sup>, Raphaël M. Kalengayi<sup>6</sup>, Eric Van Marck<sup>8</sup>, M. Thomas P. Gilbert<sup>1†</sup> & Steven M. Wolinsky<sup>4</sup>

Human immunodeficiency virus type 1 (HIV-1) sequences that pre-date the recognition of AIDS are critical to defining the time of origin and the timescale of virus evolution<sup>1,2</sup>. A viral sequence from 1959 (ZR59) is the oldest known HIV-1 infection<sup>1</sup>. Other historically documented sequences, important calibration points to convert evolutionary distance into time, are lacking, however; ZR59 is the only one sampled before 1976. Here we report the amplification and characterization of viral sequences from a Bouin's-fixed paraffin-embedded lymph node biopsy specimen obtained in 1960 from an adult female in Léopoldville, Belgian Congo (now Kinshasa, Democratic Republic of the Congo (DRC)), and we use them to conduct the first comparative evolutionary genetic study of early pre-AIDS epidemic HIV-1 group M viruses. Phylogenetic analyses position this viral sequence (DRC60) closest to the ancestral node of subtype A (excluding A2). Relaxed molecular clock analyses incorporating DRC60 and ZR59 date the most recent common ancestor of the M group to near the beginning of the twentieth century. The sizeable genetic distance between DRC60 and ZR59 directly demonstrates that diversification of HIV-1 in west-central Africa occurred long before the recognized AIDS pandemic. The recovery of viral gene sequences from decades-old paraffin-embedded tissues opens the door to a detailed palaeovirological investigation of the evolutionary history of HIV-1 that is not accessible by other methods.

We screened 27 tissue blocks (8 lymph node, 9 liver and 10 placenta) obtained from Kinshasa between 1958 and 1960 by polymerase chain reaction with reverse transcription (RT-PCR); one lymph node biopsy specimen contained HIV-1 RNA. Viral nucleic acids were extracted from this specimen using protocols optimized for the recovery of nucleic acids from ancient or degraded samples<sup>3,4</sup>. After reverse transcription, 12 out of the 14 short HIV-1 complementary DNA fragments in the study (Fig. 1a) were amplified by PCR using a panel of conserved primer pairs from different regions of the viral genome (Supplementary Table 1). Each PCR product was cloned and sequenced. Sequences were reproducible after repeated extractions and were not the result of PCR contamination (see Fig. 1a and Supplementary Table 1 for fragment designations). The results were confirmed independently in two laboratories (Fig. 1b and Supplementary Fig. 1), with the second laboratory successfully identifying the positive 1960 specimen in a blinded assay. The short fragments of the 1960 sample were found to be of subtype A and not to be a mosaic of contemporary sequences (see Supplementary Information for a detailed discussion of the authenticity of the 1960



**Figure 1 | Fragments amplified from DRC60, and the results of the phylogenetic and sequence analyses.** **a**, The HIV-1 genome fragments that were successfully amplified from DRC60 (red) and are available for ZR59 (black). The numbering for the HIV-1 sequences corresponds to the HXB2 reference sequence (Supplementary Table 1). **b**, The A/A1 subtree from the unconstrained (in which a molecular clock is not enforced) BMCMC phylogenetic analysis. Supplementary Fig. 1 depicts the complete phylogenetic tree (50% majority rule consensus tree of the posterior sample, with branch lengths averaged across the sample). Posterior probabilities are shown on nodes with support >0.95. 1960.DRC60A is the University of Arizona consensus sequence, and 1960.DRC60N is the Northwestern University consensus sequence (that is, the sequences independently recovered in each of the two laboratories). The DRC60 sequences form a strongly supported clade with three modern sequences also sampled in the DRC. **c**, Smoothed histograms of within-subtype (A2, A/A1, B, C, D, F1, F2, H, J, K) and between-subtype distances.

<sup>1</sup>Ecology and Evolutionary Biology, University of Arizona, Tucson, Arizona 85721, USA. <sup>2</sup>Sanofi Pasteur, F-69367 Lyon Cedex 07, France. <sup>3</sup>UCB SA Pharma, Braine l'Alleud, BE-1420, Belgium. <sup>4</sup>The Feinberg School of Medicine, Northwestern University, Chicago, Illinois 60611, USA. <sup>5</sup>Ancient DNA Laboratory, School of Biological Sciences and Biotechnology, Murdoch University, Perth, Western Australia 6150, Australia. <sup>6</sup>Department of Anatomy and Pathology, University of Kinshasa, Kinshasa B.P. 864, Democratic Republic of the Congo. <sup>7</sup>National Institute for Biomedical Research, National Laboratory of Public Health, Kinshasa B.P. 1197, Democratic Republic of the Congo. <sup>8</sup>Department of Pathology, University Hospital, University of Antwerp, Antwerp B-2610, Belgium. †Present address: Centre for Ancient Genetics, Biological Institute, University of Copenhagen, Copenhagen DK-2100, Denmark.

sequences). Consensus nucleotide sequences from these short HIV-1 fragments were concatenated for study. The analyses included reference sequences from the Los Alamos National Laboratory HIV sequence database and sequences recovered as part of this study from three paraffin-embedded tissue specimens collected from AIDS patients in Belgium and Canada between 1981 and 1997.

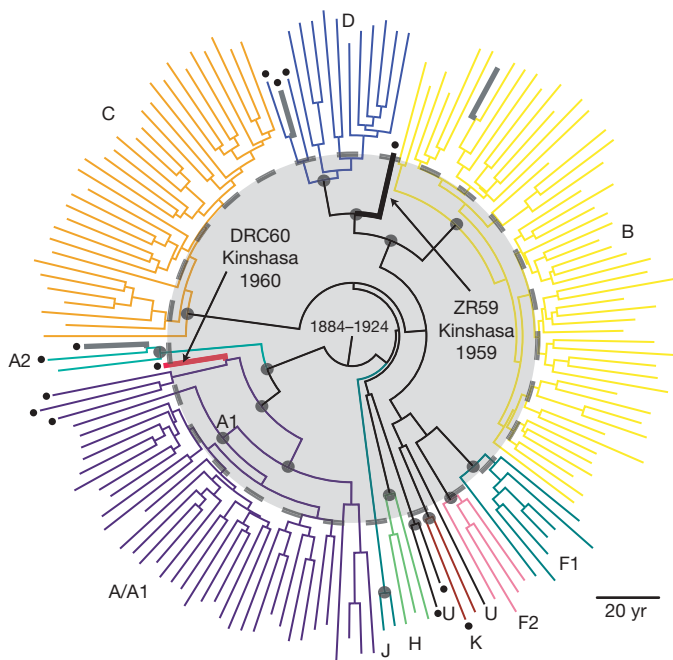
HIV-1 sequences were analysed in MrBayes v3.1.2 (ref. 5) using an unconstrained (in which a molecular clock was not enforced) Bayesian Markov chain Monte Carlo (BMCMC) method. The phylogenetic analyses confirmed that the DRC60 consensus sequences from the two laboratories were derived from a single patient (uncorrected pairwise distance of 1.4%). The sequences were positioned close to the ancestral node of the subtype A lineage (excluding sub-subtype A2), forming a monophyletic clade with three modern sequences from the DRC (Fig. 1b and Supplementary Fig. 1). Assuming a similar rate of evolution along all branches on a tree, the divergence between two sequences reflects the time elapsed since their shared ancestor. As predicted, the DRC60 sequences had a shorter branch length to the A/A1 ancestral node than the contemporary subtype A viruses sampled from the same geographic region ( $P = 1.0$ ).

We validated the time of origin of the 1960 sequence by comparisons of the predicted date to the documented date. With the DRC60 date treated as an unknown, we calculated an evolutionary rate on the basis of the distribution of branch lengths on the unconstrained phylogenetic trees sampled by MrBayes. To limit the effects of evolutionary rate differences between clades and uncertainties in rooting the HIV-1 M group phylogeny, we focused on the subtype A/A1 subtree (Supplementary Fig. 1) and analysed root-to-tip branch lengths relative to the sampling year. The mean estimates for the year of origin of the DRC60 consensus sequences from the University of Arizona and Northwestern University laboratories were 1959 (95% highest probability distribution (HPD) 1902–1984) and 1959 (95% HPD 1915–1985), respectively, corroborating the authenticity of the DRC60 sequences and the existence of a clock-like signal in our data set (see later). Despite initial indications that recombination might seriously confound phylogenetic dating estimates<sup>6</sup>, subsequent work has suggested that recombination is not likely to systematically bias HIV-1 dates in one direction or the other, although it is expected to increase variance<sup>7</sup>. The close match between the predicted and the actual dates of both ZR59 (ref. 2) and DRC60 provides support for this view and gives an unambiguous indication that HIV-1 evolves in a fairly reliable clock-like fashion.

The uncorrected pairwise distance between DRC60 and ZR59 in their overlapping *env* region was 11.7% (Fig. 1c). This genetic distance is greater than 99.2% of within-subtype comparisons (within-subtype difference, range 0.01–0.15; between-subtype difference, range 0.05–0.18). Because each subtype represents several decades of independent evolution in the human population<sup>2,8</sup>, the extensive divergence between DRC60 and ZR59 indicates that the HIV-1 M group founder virus began to diversify in the human population (and that HIV-1 probably entered Kinshasa) decades before 1960.

We applied a relaxed clock BMCMC coalescent framework as implemented in BEAST v1.4.7 (ref. 9) to estimate the time to the most recent common ancestor (TMRCA) of the HIV-1 M group. This approach robustly incorporates phylogenetic uncertainty and accounts for the possibility of variable substitution rates among lineages and differences in the demographic history of the virus, sampling phylogenies and parameter estimates in proportion to their posterior probability<sup>10</sup>. As with other studies of HIV-1 (ref. 11), comparisons of the marginal likelihoods of strict versus relaxed clock models (both of which are implemented in BEAST) indicated overwhelming support for relaxed clocks (data available on request). Hence, the use of strict clock models with these data would be inappropriate and would probably yield misleadingly small error estimates with regard to both timing and substitution rates.

Using substitution rates calibrated with sequences sampled at different time points, we obtained a posterior distribution of rooted tree topologies with branch lengths in unit time (Fig. 2 and Supplementary Fig. 2). The median estimated substitution rate for the concatenated subregions of the *gag-pol-env* genes was  $2.47 \times 10^{-3}$  substitutions per site per year (95% HPD  $1.90$ – $2.95 \times 10^{-3}$ ). The inclusion of the 1959 and 1960 sequences seemed to improve estimation of the TMRCA of the M group (Table 1), limiting the influence of the coalescent tree prior on the posterior TMRCA distributions compared with the data set that excluded these earliest cases of HIV-1. With DRC60 and ZR59 included, the different demographic/coalescent models gave highly consistent results, with tighter and more similar date ranges compared with the analyses that excluded them and 95% HPDs that extend no later than 1933. The best-fit model incorporated a constant population size demographic model (TMRCA 1921, 95% HPD 1908–1933). The model with a general, non-parametric prior (the Bayesian skyline plot tree prior)<sup>12,13</sup> that indicated a more complex (and biologically plausible) demographic history (Supplementary Fig. 3) had a statistically indistinguishable degree of support (TMRCA 1908, 95% HPD 1884–1924). Moreover, the population expansion demographic model<sup>9</sup>, which was a slightly worse fit to the data compared with the constant population and Bayesian skyline plot models, could not be rejected given the Bayes factor comparison of models (Table 1). The inability to strongly reject the model with a constant population size prior is counterintuitive because it is clear that the HIV-1 population size has increased notably. We speculate that this finding might be due to the simplest model providing a good fit to a relatively short,



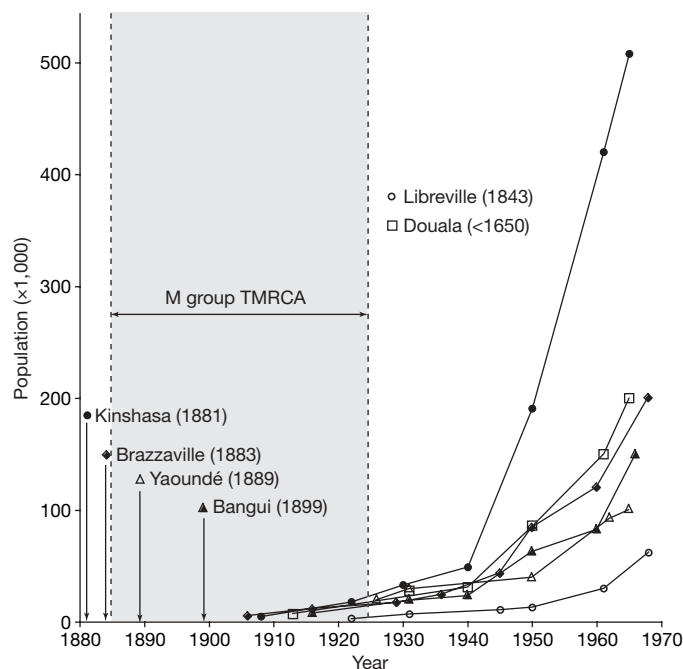
**Figure 2 | Maximum clade credibility topology inferred using BEAST v1.4.7 under a Bayesian skyline plot tree prior.** Branch lengths are depicted in unit time (years) and represent the median of those nodes that were present in at least 50% of the sampled trees. DRC60 (red), ZR59 (black) and the three control sequences from paraffin-embedded specimens from known AIDS patients (grey) are depicted in bold. The 95% HPD of the TMRCA is indicated at the root of the tree. Nodes (sub-subtype and deeper) with posterior probability of 1.0 are marked with grey circles. Unclassifiable strains are labelled 'U'. Sequences sampled in the DRC are highlighted with a bullet at the tip. DRC60 and the two control sequences from the DRC each form monophyletic clades with previously published sequences from the DRC, whereas the Canadian control sequence clusters, as expected, with subtype B sequences. The dashed circle and shaded area show the extensive HIV-1 diversity in Kinshasa in the 1950s. Supplementary Fig. 2 shows the tree in rectangular form with taxon labels.

information-poor alignment, in comparison with more parameterized models.

Acid-containing fixatives such as Bouin's solution can cause base modifications of nucleic acids, leading to the generation of erroneous bases in sequences derived from such samples<sup>3</sup>. However, the replication of all sequences from independent PCR amplifications and the uncorrected pairwise distance between the consensus sequences from the two laboratories (1.4%) suggest that few of the mutations on the DRC60 lineage are damaged-induced. Moreover, our relaxed clock methods are likely to be fairly robust to the presence of such mutations in one lineage<sup>9</sup>. Nevertheless, additional old sequence data would be helpful for resolving what impact, if any, this possible source of error had on the slightly earlier dates we calculated compared with previous estimates that did not include early calibration points<sup>2,8,14,15</sup>. Interestingly, the best-fit model for the data set that excluded ZR59 and DRC60 (Table 1) gave a TMRCA estimate of 1933 (1919–1945), which is very similar to that of ref. 2. This suggests that the inclusion of the old sequences, rather than the vagaries associated with a much shorter alignment than that analysed by ref. 2, might explain the discrepancy. Also, one earlier study, using sequences from the DRC only<sup>16</sup>, produced dating and demography estimates very similar to ours. Overall, there is broad agreement between all of these studies in spite of differences in data and methods.

Our estimation of divergence times, with an evolutionary timescale spanning several decades, together with the extensive genetic distance between DRC60 and ZR59 indicate that these viruses evolved from a common ancestor circulating in the African population near the beginning of the twentieth century; TMRCA dates later than the 1930s are strongly rejected by our statistical analyses. The topology of the HIV-1 group M phylogeny provides further support for this conclusion. Unlike ZR59, which is basal to subtype D<sup>1</sup>, DRC60 branches off from the ancestral node of subtype A/A1 (Fig. 2 and Supplementary Figs 1 and 2). Thus, it is clear that phylogenetically distinct subtypes (and/or their progenitors) were already present in the DRC by this early time point (Fig. 2). Notably, DRC60 and ZR59 cluster with other strains from the same geographical region and basal to other members of their respective subtypes, a pattern consistent with the hypothesis that the subtypes spread through lineage founder effects worldwide, whereas a more diverse array of forms remained at the site of origin in Africa<sup>17,18</sup>.

The reservoir of the ancestral virus still exists among wild chimpanzee communities in the same area on the African continent<sup>19</sup>. Humans acquired a common ancestor of the HIV-1 M group by cross-species transmission under natural circumstances<sup>20</sup>, probably predation<sup>21</sup>. The Bayesian skyline plot (Supplementary Fig. 2), which tracks effective population size through time, suggests that HIV-1 group M experienced an extensive period of relatively slow growth in the first half of the twentieth century. A similar pattern has been inferred using sequences sampled only in the DRC<sup>16</sup>. This pattern, and the short duration between the first presence of urban agglomerations in this area and the timing of the most recent common ancestor of HIV-1 group M (Fig. 3), suggests that the rise of cities may have facilitated the initial establishment and the early spread of HIV-1. Hence, the founding and growth of colonial administrative and trading centres such as Kinshasa<sup>22</sup> may have enabled the region to become the epicentre of the HIV/AIDS pandemic<sup>23</sup>.



**Figure 3 | The origin and growth of the major settlements near the epicentre of the HIV-1 group M epidemic.** In the countries surrounding the putative zone of cross-species transmission<sup>19</sup> (current-day Cameroon, Central African Republic, DRC, Republic of Congo, Gabon and Equatorial Guinea) there was not a single site with a population exceeding 10,000 until after 1910. The founding date of each major city in the region is listed beside its name. Most were founded only shortly before the estimated TMRCA of group M. The demographic data are from ref. 23.

The archival banks of Bouin's-fixed paraffin-embedded tissue specimens accumulated by many hospitals in west-central Africa provide a vast source of clinical material for viral genetic analysis. As with the 1918 Spanish influenza pandemic virus<sup>24,25</sup>, a deep perspective on the evolutionary history of HIV-1 using sequences resurrected from the earliest cases in Africa could yield important insights into the pathogenesis, virulence and evolution of pandemic AIDS viruses.

## METHODS SUMMARY

A total of 813 Bouin's-fixed paraffin-embedded histopathological blocks were recovered from the 1958–1962 archives of the Department of Anatomy and Pathology at the University of Kinshasa. The boxes were stored until transfer to the University of Arizona, where 8 lymph node, 9 liver and 10 placenta samples from 1958–1960 were selected for RNA preservation analysis and HIV-1 RNA screening. We used a human  $\beta$ -2-microglobulin (*B2M*) quantitative RT-PCR assay to assess RNA quality as described<sup>3</sup>. Digestion and extraction of these samples, and of three modern positive-control samples, were performed using QIAamp DNA micro kits (Qiagen) using the protocol described in ref. 3. We used 14 primer sets designed to anneal to highly conserved regions of the *gag*, *pol* and *env* genes of HIV-1 group M and to amplify very short fragments likely to be present even in ancient and/or degraded specimens (Supplementary Table 1). Reverse transcription was performed using the SuperScript III System for RT-PCR (Invitrogen). The cDNA was amplified by PCR using Platinum Taq HiFi enzyme (Invitrogen) and cloned using the TOPO TA Cloning Kit (Invitrogen). We constructed an alignment including 156 published reference sequences plus

**Table 1 | HIV-1 M group TMRCA estimates from BEAST analyses under different coalescent tree priors**

Coalescent tree prior	DRC60 and ZR59 excluded*	DRC60 and ZR59 included
Constant	<b>1933 (1919–1945)†, 0.0</b>	<b>1921 (1908–1933)†, 0.0</b>
Exponential	1907 (1874–1932), $-3.5 \pm 0.8$	1914 (1891–1930), $-2.1 \pm 1.5$
Expansion	1882 (1834–1917), $-2.7 \pm 0.8$	<b>1902 (1873–1922)†, <math>-1.6 \pm 1.5</math></b>
Logistic	1913 (1880–1937), $-2.3 \pm 0.8$	1913 (1891–1930), $-3.2 \pm 1.5$
Bayesian skyline plot	1882 (1831–1916), $-2.7 \pm 0.8$	<b>1908 (1884–1924)†, <math>-0.4 \pm 1.5</math></b>

Shown for each coalescent tree prior is the median, with the 95% highest probability distribution of TMRCA in parentheses. Also shown is the  $\log_{10}$  Bayes factor difference in estimated marginal likelihood ( $\pm$  estimated standard error) compared with the coalescent model with strongest support.

\*Concatenated *gag-pol-env* fragments available for either or both of ZR59 and DRC60 (994 nucleotides total, 507 from DRC60).

†TMRCA for the best-fit model and models not significantly worse than it are written in bold.



the sequences recovered in this study, concatenating the 12 (out of 14) fragments successfully amplified from the 1960 sample and the 4 fragments already available from the 1959 sample (994 bases total). We performed an unconstrained (not enforced by a molecular clock) BMCMC analysis in MrBayes v3.1.2 (ref. 5) and used the resulting MCMC sample to test whether the 1960 sequence exhibited properties consistent with its provenance (both age and geography). We used a relaxed molecular clock model, as implemented in BEAST v1.4.7 (ref. 9), to estimate the TMRCA of HIV-1 group M using the 1960 and 1959 samples and to investigate the demographic history of the virus. We also performed pairwise comparisons within and between subtypes for the 163 bases available for both DRC60 and ZR59.

**Full Methods** and any associated references are available in the online version of the paper at [www.nature.com/nature](http://www.nature.com/nature).

**Received 21 May; accepted 8 September 2008.**

- Zhu, T. F. *et al.* An African HIV-1 sequence from 1959 and implications for the origin of the epidemic. *Nature* **391**, 594–597 (1998).
- Korber, B. *et al.* Timing the ancestor of the HIV-1 pandemic strains. *Science* **288**, 1789–1796 (2000).
- Gilbert, M. T. P. *et al.* The isolation of nucleic acids from fixed, paraffin-embedded tissues — which methods are useful when? *PLoS ONE* **2**, e537 (2007).
- Worobey, M. Phylogenetic evidence against evolutionary stasis and natural abiotic reservoirs of influenza A virus. *J. Virol.* **82**, 3769–3774 (2008).
- Huelsenbeck, J. P. & Ronquist, F. MrBayes: Bayesian inference of phylogeny. *Bioinformatics* **17**, 754–755 (2001).
- Worobey, M. A novel approach to detecting and measuring recombination: insights into evolution in viruses, bacteria, and mitochondria. *Mol. Biol. Evol.* **18**, 1425–1434 (2001).
- Lemey, P. *et al.* The molecular population genetics of HIV-1 group O. *Genetics* **167**, 1059–1068 (2004).
- Gilbert, M. T. P. *et al.* The emergence of HIV-1 in the Americas and beyond. *Proc. Natl Acad. Sci. USA* **104**, 18566–18570 (2007).
- Drummond, A. J. & Rambaut, A. BEAST: Bayesian evolutionary analysis by sampling trees. *BMC Evol. Biol.* **7**, 214 (2007).
- Drummond, A. J., Ho, S. Y. W., Phillips, M. J. & Rambaut, A. Relaxed phylogenetics and dating with confidence. *PLoS Biol.* **4**, e88 (2006).
- Salemi, M., de Oliveira, T., Ciccozzi, M., Rezza, G. & Goodenow, M. M. High-resolution molecular epidemiology and evolutionary history of HIV-1 subtypes in Albania. *PLoS ONE* **3**, e1390 (2008).
- Suchard, M. A., Weiss, R. E. & Sinsheimer, J. S. Bayesian selection of continuous-time Markov chain evolutionary models. *Mol. Biol. Evol.* **18**, 1001–1013 (2001).
- Drummond, A. J., Rambaut, A., Shapiro, B. & Pybus, O. G. Bayesian coalescent inference of past population dynamics from molecular sequences. *Mol. Biol. Evol.* **22**, 1185–1192 (2005).
- Sharp, P. M. *et al.* The origins of acquired immune deficiency syndrome viruses: where and when? *Phil. Trans. R. Soc. Lond. B* **356**, 867–876 (2001).
- Salemi, M. *et al.* Dating the common ancestor of SIVcpz and HIV-1 group M and the origin of HIV-1 subtypes using a new method to uncover clock-like molecular evolution. *FASEB J.* **15**, 276–278 (2001).
- Yusim, K. *et al.* Using human immunodeficiency virus type 1 sequences to infer historical features of the acquired immune deficiency syndrome epidemic and human immunodeficiency virus evolution. *Phil. Trans. R. Soc. Lond. B* **356**, 855–866 (2001).
- Vidal, N. *et al.* Unprecedented degree of human immunodeficiency virus type 1 (HIV-1) group M genetic diversity in the Democratic Republic of Congo suggests that the HIV-1 pandemic originated in Central Africa. *J. Virol.* **74**, 10498–10507 (2000).
- Rambaut, A., Robertson, D. L., Pybus, O. G., Peeters, M. & Holmes, E. C. Human immunodeficiency virus phylogeny and the origin of HIV-1. *Nature* **410**, 1047–1048 (2001).
- Keele, B. F. *et al.* Chimpanzee reservoirs of pandemic and nonpandemic HIV-1. *Science* **313**, 523–526 (2006).
- Worobey, M. in *Global HIV/AIDS Medicine* (eds Volberding, P. A., Sande, M. A., Lange, J. & Greene, W. C.) 13–21 (Saunders Elsevier, 2008).
- Hahn, B. H., Shaw, G. M., De Cock, K. M. & Sharp, P. M. AIDS as a zoonosis: scientific and public health implications. *Science* **287**, 607–614 (2000).
- Hance, W. A. *Population, Migration, and Urbanization in Africa* 209–297 (Columbia Univ. Press, 1970).
- Chitnis, A., Rawls, D. & Moore, J. Origin of HIV type 1 in colonial French Equatorial Africa? *AIDS Res. Hum. Retrov.* **16**, 5–8 (2000).
- Taubenberger, J. K. *et al.* Characterization of the 1918 influenza virus polymerase genes. *Nature* **437**, 889–893 (2005).
- Tumpey, T. M. *et al.* Characterization of the reconstructed 1918 Spanish Influenza pandemic virus. *Science* **310**, 77–80 (2005).

**Supplementary Information** is linked to the online version of the paper at [www.nature.com/nature](http://www.nature.com/nature).

**Acknowledgements** We thank J. Wertheim and M. Sanderson for computational assistance, and L. Jewel for providing the Canadian control specimen. The NIH/NIAD and the David and Lucile Packard Foundation funded the research.

**Author Contributions** M.W., D.E.T., S.M.W. and M.T.P.G. designed the study. M.G., T.H., K.K. and M.T.P.G. performed digestion and extraction, PCR, quantitative PCR, cloning and sequencing experiments. M.T.P.G., M.G. and M.B. optimized DNA/RNA isolation methods and designed PCR assays. D.E.T., J.-J.M., E.V.M., J.-M.M.K. and R.M.K. organized and provided samples. M.W. analysed the data, performed the phylogenetic analyses, and wrote the paper. S.M.W. contributed to the analyses and writing. All authors discussed the results and commented on the manuscript.

**Author Information** The sequences reported in this study have been deposited in GenBank under accession numbers EU580739–EU580854 and EU589211–EU589236. Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). Correspondence and requests for materials should be addressed to M.W. ([worobey@email.arizona.edu](mailto:worobey@email.arizona.edu)).

## METHODS

**Archival samples.** Each individual block carried an original paper identification number permanently embedded in the paraffin. Laboratory books listed the corresponding identification numbers sequentially, and included the patient's age, sex, department of hospitalization, tissue type and date of sampling. Block identification number, sampling date and tissue type were transcribed onto an Excel spreadsheet, and the blocks were indexed, transferred into plastic boxes and photographed.

The results of the quantitative RT-PCR assay indicated that the integrity of the RNA preserved in these 27 samples ranged from moderate to undetectable, a range typical of Bouin's-fixed specimens<sup>3</sup>. The human RNA found in the 1960 lymph node biopsy sample that was found to be HIV-1-RNA-positive was of relatively good quality. The  $C_t$  values (quantitative PCR data available from the authors on request) were as low or lower (better) than more recent (1980–1990) paraffin-embedded tissues that have yielded short HIV-1 RNA amplicons<sup>3</sup>.

Three formalin-fixed paraffin-embedded necropsy specimens were obtained: a Canadian patient who died in 1997 (CAN97); a Congolese woman who died in Belgium in 1981 and who was retrospectively identified as an AIDS patient (BE81); and a Congolese man who died in Belgium in 1985 (BE85). The latter two cases were presumably infected in Zaire (now the Democratic Republic of the Congo). The phylogenetic reconstruction shows that their viral sequences are most closely related to modern sequences from the Democratic Republic of the Congo, whereas the Canadian specimen yielded a subtype B sequence, as predicted (Figs 1 and 2 and Supplementary Figs 1 and 2).

**RNA isolation and reverse transcription.** Between 5 and 10 microtome sections, 5–10 µm in thickness, or an approximately equivalent amount of tissue shaved from each block with a disposable scalpel blade, were used for each digestion and extraction, as described<sup>3</sup>. Rigorous attention was given to preventing cross-contamination between samples by cleaning the outer surface of each block with a bleach solution, using fresh microtome/scalpel blades for each sectioning of each block, discarding the first few (exposed-surface) sections, and by performing the work in a room physically isolated from any human or HIV-1 PCR-product DNA. A 48-h digestion period (24 h at 65 °C, 24 h at 75 °C) was used. Post extraction nucleic acids were eluted into 100 µl elution buffer AE and stored frozen at –80 °C until required for analyses.

Reverse transcription was performed simultaneously for the *gag*, *pol* and human B2M RNA fragments; *env* fragments 3–10; and *env* fragments 11–14 (Supplementary Table 1). This was performed with SuperScript III used according to the manufacturer's instructions. The protocol was as described<sup>3</sup> except that alternating 50 °C and 55 °C incubation periods of 30 min were used for a total of 6 h.

**Amplification, cloning and DNA sequencing.** The cDNA was PCR amplified in 25-µl reactions, using 0.1 µl Platinum Taq HiFi enzyme (Invitrogen), 250 µM dNTP mix, 2 mM MgSO<sub>4</sub>, 1× PCR buffer, 0.4 µM per primer, and 2 µl cDNA for the *gag* and *pol* reactions or 1 µl for the *env* ones, with annealing temperatures of 60 °C (*gag*, 60 cycles) or 55 °C (*pol*, 50 cycles; *env*, 55 cycles). Full details are available from the authors on request.

After amplification, the PCR-product DNA was visualized by agarose gel electrophoresis and then purified using Zymoclean DNA Clean and Concentrator-25 spin tubes (Zymo Research Corporation). PCR-product DNA was inserted into vector pCR2.1-TOPO using the TOPO TA Cloning Kit (Invitrogen). The University of Arizona Genomic Analysis and Technology Core Facility resolved the DNA sequence of the vector inserts on an Applied Biosystems 3730xl DNA analyser using ABI Big Dye 3.1 chemistry (Applied Biosystems). Nearly identical protocols were followed for the independent replication of the DRC60 results at Northwestern University.

**Alignments.** We downloaded the 2006 full-length HIV-1 sequence alignment from the Los Alamos National Laboratories HIV sequence database<sup>26</sup>. We retained only non-recombinant HIV-1 group M A–K subtype sequences (excluding G) and removed sequences suspected a priori of unusual evolutionary dynamics (such as those associated with the intravenous drug user epidemic in Eastern Europe and those with *nef* deletions, both of which exhibit abnormally slow evolutionary rates). We also reduced the size of the subtype B and C clades, which are heavily over-sampled relative to the others, by keeping only the first 5 sequences from any year/country pair and then randomly removing sequences until the sample size was similar to that of the other subtypes. This procedure left a total of 156 sequences. We then manually aligned the consensus sequence from

the 12 regions amplified from DRC60, plus the 4 regions available for ZR59, to the full-length sequences. These short regions (Fig. 1a and Supplementary Table 1) were then concatenated into an alignment 994 nucleotides in length. The four *env* fragments from DRC60 that overlapped with available data from ZR59 were concatenated into an alignment 163 nucleotides in length. Matching alignments with DRC60 and ZR59 removed were also constructed. All the alignments are available from the authors on request.

**MrBayes analyses.** We used a general time-reversible nucleotide substitution model with gamma-distributed rate heterogeneity among sites and performed four independent runs of 20 million steps, sampling every 2,000 steps. Examination of the MCMC samples with Tracer v1.4 (ref. 9) indicated convergence and adequate mixing of the Markov chain with estimated sample sizes in the thousands. We discarded the first 2 million steps from each run as burn-in, and combined the resulting MCMC samples for subsequent estimation of posteriors. The 50% majority rule consensus tree (Supplementary Fig. 1) is shown rooted on the branch identified by the rooted-tree method in BEAST v1.4.7 (ref. 9), described below; however, the group M rooting was not relevant to any dating analysis. We also estimated phylogenies using the same data set under neighbour-joining and maximum likelihood methods and the same substitution model. The DRC60 sequences fell in the same topological position as with the BMCMC methods, with short root-to-tip genetic distances, consistent with the MrBayes results (Supplementary Fig. 1). All data and trees available from the authors on request.

We used the posterior tree sample to test the hypothesis that the terminal nodes of the DRC60 sequences were closer to the inferred A/A1 ancestral node by calculating the proportion of sampled trees where the A/A1 node-to-tip distances were smaller for these sequences than for the three modern sequences from the DRC in the same clade (Fig. 1b).

To predict the date of sampling on the basis of the phylogenetic properties of the DRC60 sequences, we also plotted the branch lengths (A/A1 node to tips) against the time of sampling for all A/A1 sequences excluding DRC60 and calculated the best fit for the linear regression of genetic divergence against the year of sampling of the viruses<sup>27</sup>. We calculated the mean and 95% HPD of the predicted sampling date of each DRC60 consensus sequence on the basis of its node-to-tip distance and the inferred regression line calculated for each of 100 trees sampled by MrBayes.

**Bayesian MCMC inference of phylogeny using BEAST v1.4.7.** We used the Bayesian methods described previously<sup>9,10</sup>, which allow for the co-estimation of phylogeny and divergence times under a 'relaxed' molecular clock model, as implemented in BEAST v1.4.7 (ref. 9). All analyses were performed under an uncorrelated lognormal relaxed molecular clock model, using a general time-reversible nucleotide substitution model with heterogeneity among sites modelled with a gamma distribution. We investigated each demographic model (constant population, exponential growth, expansion growth, logistic growth) as well as a Bayesian skyline plot coalescent tree prior<sup>13</sup>, a general, non-parametric prior that enforces no particular demographic history. We used a piecewise linear skyline model with 10 groups. We then compared the marginal likelihoods for each model using Bayes factors estimated in Tracer v1.4 as described<sup>12,15</sup>. Bayes factors represent the ratio of the marginal likelihoods of the models being compared. A large ratio can indicate that one model is a significantly better fit to the data than another. We assessed the strength of the evidence that the best-fit model was superior to the others as described<sup>15</sup>.

For each analysis, two independent runs of 50 million steps were performed. Examination of the MCMC samples with Tracer v1.4 indicated convergence and adequate mixing of the Markov chains, with estimated sample sizes in the hundreds or thousands. After inspection with Tracer, we discarded an appropriate number of steps from each run as burn-in, and combined the resulting MCMC tree samples for subsequent estimation of posteriors. We summarized the MCMC samples using the maximum clade credibility topology found with TreeAnnotator v1.4.7 (ref. 9), with branch length depicted in years (median of those branches that were present in at least 50% of the sampled trees; Fig. 2). The Bayesian skyline plot was reconstructed using the posterior tree sample and Tracer v1.4.

26. Leitner, T. *et al.* HIV Sequence Compendium (<http://www.hiv.lanl.gov>) (Theoretical Biology and Biophysics Group, Los Alamos National Laboratory, 2005).

27. Drummond, A., Pybus, O. G. & Rambaut, A. Inference of viral evolutionary rates from molecular sequences. *Adv. Parasitol.* **54**, 331–358 (2003).