



## Ancient hepatitis B viruses from the Bronze Age to the Medieval period

Mühlemann, Barbara; Jones, Terry C.; Damgaard, Peter de Barros; Allentoft, Morten E.; Shevnina, Irina; Logvin, Andrey; Usmanova, Emma; Panyushkina, Irina P.; Boldgiv, Bazartseren; Bazartseren, Tsevel

*Published in:*  
Nature

*Link to article, DOI:*  
[10.1038/s41586-018-0097-z](https://doi.org/10.1038/s41586-018-0097-z)

*Publication date:*  
2018

*Document Version*  
Peer reviewed version

[Link back to DTU Orbit](#)

*Citation (APA):*  
Mühlemann, B., Jones, T. C., Damgaard, P. D. B., Allentoft, M. E., Shevnina, I., Logvin, A., ... Willerslev, E. (2018). Ancient hepatitis B viruses from the Bronze Age to the Medieval period. *Nature*, 557(7705), 418-423. <https://doi.org/10.1038/s41586-018-0097-z>

---

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

# 1 Ancient Hepatitis B viruses from the 2 Bronze Age to the Medieval

3 Barbara Mühlemann<sup>\*,1</sup>, Terry C. Jones<sup>\*,1,2</sup>, Peter de Barros Damgaard<sup>\*,3</sup>, Morten E.  
4 Allentoft<sup>\*,3</sup>, Irina Shevnina<sup>4</sup>, Andrey Logvin<sup>4</sup>, Emma Usmanova<sup>5</sup>, Irina P.  
5 Panyushkina<sup>6</sup>, Bazartseren Boldgiv<sup>7</sup>, Tsevel Bazartseren<sup>8</sup>, Kadicha Tashbaeva<sup>9</sup>,  
6 Victor Merz<sup>10</sup>, Nina Lau<sup>11</sup>, Václav Smrčka<sup>12</sup>, Dmitry Voyakin<sup>13</sup>, Egor Kitov<sup>14</sup>,  
7 Andrey Epimakhov<sup>15</sup>, Dalia Pokutta<sup>16</sup>, Magdolna Vicze<sup>17</sup>, T. Douglas Price<sup>18</sup>,  
8 Vyacheslav Moiseyev<sup>19</sup>, Anders J. Hansen<sup>3</sup>, Ludovic Orlando<sup>3,20</sup>, Simon  
9 Rasmussen<sup>21</sup>, Martin Sikora<sup>3</sup>, Lasse Vinner<sup>3</sup>, Albert D. M. E. Osterhaus<sup>22</sup>, Derek J.  
10 Smith<sup>1</sup>, Dieter Glebe<sup>23,24</sup>, Ron A. M. Fouchier<sup>25</sup>, Christian Drosten<sup>2,26</sup>, Karl-Göran  
11 Sjögren<sup>18</sup>, Kristian Kristiansen<sup>18</sup>, Eske Willerslev<sup>§,3,27,28</sup>.

12

13 1. Center for Pathogen Evolution, Department of Zoology, University of Cambridge, Downing St.,  
14 Cambridge CB2 3EJ, UK. 2. Institute of Virology, Charité, Universitätsmedizin Berlin, Campus  
15 Charité Mitte, Charitéplatz 1, 10117 Berlin, Germany. 3. Centre for GeoGenetics, Natural History  
16 Museum, University of Copenhagen, Øster Voldgade 5-7, 1350 Copenhagen K, Denmark. 4.  
17 Archaeological Laboratory, Faculty of History and Law, A. A. Baitursynov Kostanay State University,  
18 47 Baitursynov St., Kostanay, 110000, Kazakhstan 5. Karaganda State University, Saryarka  
19 Archaeological Institute, 28 Universitetskaya St., 100012 Karaganda, 100012, Kazakhstan. 6.  
20 Laboratory of Tree-Ring Research, University of Arizona, 1215 E. Lowell St., Tucson, AZ 85721,  
21 USA. 7. Department of Biology, School of Arts and Sciences, National University of Mongolia,  
22 Ulaanbaatar 14201, Mongolia. 8. Laboratory of Virology, Institute of Veterinary Medicine, Mongolian  
23 University of Life Sciences, Ulaanbaatar 17024, Mongolia. 9. National Academy of Sciences, 265a  
24 Chuy Ave., Bishkek, 720001, Kyrgyzstan. 10. Pavlodar State University, 64 Lomov Str., Pavlodar,  
25 140008, Kazakhstan. 11. Centre for Baltic and Scandinavian Archaeology, Schloss Gottorf, D 24837  
26 Schleswig, Germany. 12. Institute for History of Medicine and Foreign Languages of the First Faculty  
27 of Medicine, Charles University, 32 Kateřinská, Prague, 121 08, Czech Republic. 13. Margulan  
28 Institute of Archaeology, 71 Zenkov St., Almaty, 050010, Kazakhstan. 14. N. N. Miklouho-Maklay  
29 Institute of Ethnology and Anthropology, Russian Academy of Sciences, 32a Leninsky Ave., Moscow,  
30 1199111, Russia. 15. South Ural Department, Institute of History and Archaeology UBRAS, South  
31 Ural State University, 76 Lenina St., Chelyabinsk, 454080, Russia. 16. Department of Archaeology and  
32 Classical Studies, Stockholm University, Universitetsvägen 10, Stockholm, 11418, Sweden. 17.  
33 Matrica Museum, 1-3 Gesztenyés St., Százhalombatta, 2440, Hungary. 18. Department of Historical  
34 Studies, University of Gothenburg, Eklandagatan 86, 412 61 Göteborg, Sweden. 19. Department of

35 Physical Anthropology, Peter the Great Museum of Anthropology and Ethnography, 2 Universitetskaya  
36 Naberazhnaya, Saint-Petersburg, 199034, Russia. 20. Laboratoire d'Anthropobiologie Moléculaire et  
37 d'Imagerie de Synthèse, CNRS UMR 5288, Université de Toulouse, Université Paul Sabatier, CNRS  
38 UMR 5288, Toulouse, 31000, France. 21. Department of Bio and Health Informatics, Technical  
39 University of Denmark, 2800 Kgs Lyngby, Denmark. 22. Research Center for Emerging Infections and  
40 Zoonoses, University of Veterinary Medicine Hannover, Bünteweg 17, Hannover, 30559, Germany.  
41 23. Institute of Medical Virology, Justus Liebig University of Giessen, Schubertstrasse 81, Giessen, D  
42 35392, Germany. 24. National Reference Centre for Hepatitis B and D Viruses, German Center for  
43 Infection Research (DZIF), Schubertstrasse 81, D 35392 Giessen, Germany. 25. Department of  
44 Viroscience, Erasmus Medical Centre, Wytemaweg 80, Rotterdam, 3015 CN, Netherlands. 26.  
45 German Center for Infection Research (DZIF), Inhoffenstraße 7, Braunschweig, 38124, Germany. 27.  
46 Cambridge GeoGenetics Group, Department of Zoology, Downing St., University of Cambridge,  
47 Downing Street, Cambridge CB2 3EJ, UK. 28. Wellcome Trust Sanger Institute, Hinxton, CB10 1SA,  
48 UK.

49

50 \* These authors contributed equally to this work.

51 § To whom correspondence should be addressed.

52

53

54

55 **Abstract**

56 **Hepatitis B virus (HBV) is a major cause of human hepatitis. There is**  
57 **considerable uncertainty about the timescale of its evolution and its association**  
58 **with humans. Here we present 12 full or partial ancient HBV genomes between**  
59 **~0.8-4.5 thousand years old. The ancient sequences group either within or in**  
60 **sister relationship to extant human or other ape HBV clades. Generally, the**  
61 **genome properties follow those of modern HBV. The root of the HBV tree is**  
62 **projected to between 8.6-20.9 thousand years ago (kya), and estimate a**  
63 **substitution rate between  $8.04 \times 10^{-6}$ - $1.51 \times 10^{-5}$  nucleotide substitutions per site per**  
64 **year (s/s/y). In several cases, the geographic locations of the ancient genotypes do**  
65 **not match present day distributions. Genotypes that today are typical of Africa**  
66 **and Asia, and a subgenotype from India, are shown to have an early Eurasian**  
67 **presence. The geographic and temporal patterns we observe in ancient and**  
68 **modern HBV genotypes are compatible with well-documented human**  
69 **migrations during the Bronze and Iron Ages<sup>1,2</sup>. We show evidence for the**  
70 **creation of genotype A via recombination and a long-term association of modern**  
71 **HBV genotypes with humans, including the discovery of a human genotype that**  
72 **is now extinct. Taken together, the data expose a complexity of HBV evolution**  
73 **that is not evident when considering modern sequences alone.**

74

75 HBV is transmitted perinatally or horizontally via blood or genital fluids<sup>3</sup>. The  
76 estimated global prevalence is 3.6%, ranging from 0.01% (UK) to 22.38% (South  
77 Sudan)<sup>4</sup>. In high endemicity areas, where prevalence is > 8%, 70-90% of the adult  
78 population show evidence of past or present infection<sup>5,6</sup>. The young and the  
79 immunocompromised are most likely to develop chronic HBV infection, which can

80 result in high viremia over years to decades<sup>3</sup>. Approximately 257 million people are  
81 chronically infected<sup>5</sup> and around 887,000 died in 2015 due to associated  
82 complications<sup>5</sup>.

83

84 Despite the prevalence and public health impact of HBV, its origin and evolution  
85 remain unclear<sup>7,8</sup>. Inference of HBV nucleotide substitution rates is complicated by  
86 the fact that the virus genome consists of four overlapping open reading frames<sup>9</sup>, and  
87 mutation rates differ between phases of chronic infection<sup>10</sup>. Studies based on  
88 heterochronous sequences, sampled over a relatively short time period, find higher  
89 substitution rates, whereas rates estimated using external calibrations tend to be lower,  
90 leading to a wide range of estimated substitution rates ( $7.72 \times 10^{-4}$ - $3.7 \times 10^{-6}$ ) for HBV<sup>11-</sup>  
91 <sup>13</sup>. Human HBV is classified into at least nine genotypes, A-I, roughly corresponding  
92 to sequence similarity of at least 92.5% within genotypes<sup>14</sup>, with a heterogeneous  
93 global distribution (Fig. 1a)<sup>8,9</sup>. Attempts to explain the origin of genotypes using  
94 human migrations have been inconclusive. The hypothesis that HBV co-evolved with  
95 ancient modern humans as they left Africa has been contested due to the basal  
96 phylogenetic position of genotypes F and H, found exclusively in the Americas<sup>7</sup>.

97 HBV also infects non-human primates (NHP), and the human and other great ape  
98 HBV are interspersed in the phylogenetic tree, possibly due to cross-species  
99 transmission<sup>15</sup>. Given the variability of estimated substitution rates, the incongruence  
100 of the tree topology with some human migrations, and the mixed topology of the NHP  
101 and human HBV sequences in the phylogenetic tree, considerable uncertainty remains  
102 about the evolutionary history of HBV.

103

104 Recent advances in the sequencing of ancient DNA (aDNA) have yielded important  
105 insights into human evolution, past population dynamics<sup>16</sup>, and diseases<sup>17,18</sup>.  
106 However, ancient sequences have been recovered for only a handful of exogenous  
107 human viruses, including influenza (~100 years)<sup>19</sup>, variola (~350 years)<sup>20</sup>, and HBV  
108 (~340 years and ~450 years)<sup>21,22</sup>. The knowledge gained from these few cases  
109 emphasizes the general importance of ancient sequences for the direct study of long-  
110 term viral evolution. HBV has several characteristics that make it a good candidate for  
111 detection in an aDNA virus study: its extended high viremia during chronicity<sup>3</sup>, the  
112 relative stability of its virion<sup>23</sup>, and its small, circular, and partially double-stranded  
113 DNA genome<sup>9</sup>.

114

115 Shotgun sequence data were generated from 167 Bronze Age<sup>1</sup> and 137 predominantly  
116 Iron Age<sup>2</sup> individuals from Central to Western Eurasia with a sample age range of  
117 ~7.1-0.2 kya. We identified reads that matched the HBV genome in 25 samples  
118 (Table 1, Extended Data Table 1a, SI Table 3), spanning a period of almost 4000  
119 years, from several different cultures and a broad geographical range (Fig. 1b, Table  
120 1, Extended Data Table 1a, SI Table 3). Using TaqMan PCR, we tested two samples  
121 with high genome coverage (DA195, DA222) and two samples with low coverage  
122 (DA85, DA89) for the presence of HBV. The high-coverage samples tested positive,  
123 whereas the low-coverage samples tested negative (Extended Data Table 1b). This is  
124 consistent with shotgun sequencing being more effective than targeted PCR for  
125 analysing highly degraded DNA<sup>24</sup>. Based on availability of sample material, libraries  
126 from 14 samples were selected for targeted enrichment (capture) of HBV DNA  
127 fragments (SI Tables 1 and 2). This resulted in increased genome coverage and an  
128 average of a 2.4-fold increase in number of HBV positive reads (Extended Data Table

129 1a, SI Table 3). In total, we obtained 17.9 to 100% HBV genome coverage from the  
130 sequence data, with genomic depth ranging from 0.4x to 89.2x (Table 1, Extended  
131 Data Table 1a). We selected 12 samples for phylogenetic analyses. Criteria for  
132 inclusion were at least 50% genome coverage and clear aDNA damage patterns after  
133 capture (Extended Data Fig. 1).

134

135 For an initial phylogenetic grouping, we estimated a Maximum Likelihood (ML) tree  
136 using the ancient HBV genomes together with modern human, NHP, rodent, and bat  
137 HBV genomes (Dataset 1, see Methods). All ancient viruses fell within the diversity  
138 of Old World primate HBV genotypes, which includes all human and other great ape  
139 genotypes, except human genotypes F and H (Extended Data Fig. 2).

140

141 Recombination is known to occur in HBV<sup>25</sup>. We found strong evidence that an  
142 ancient sequence, HBV-DA51, and an unknown parent recombined to form the  
143 ancient genotype A sequences. Although that cannot literally be the case due to  
144 sample ages, the logical interpretation is that an ancestor of HBV-DA51 was involved  
145 in the recombination. The same recombination is also suggested for the two modern  
146 genotype A sequences that were included in the analysis. The ancient genotype B  
147 (HBV-DA45), a modern genotype B, and two modern genotype C sequences were not  
148 similarly flagged, suggesting that the possible recombination occurred after genotypes  
149 A, B, and C had diverged. The predicted recombination break points (Extended Data  
150 Table 2, Extended Data Fig. 3) correspond closely to the polymerase gene. Thus, it is  
151 possible that the polymerase from an unknown parent and the remainder of the  
152 genome from an HBV-DA51 ancestor recombined to form the now ubiquitous  
153 genotype A about 7.4-9 kya (Fig. 2, Extended Data Table 3b, Methods). Similar

154 recombinations events, involving the creation of genotypes E and G and a currently  
155 circulating B/C recombinant, have also been identified<sup>25</sup>.

156

157 For detailed phylogenetic analyses, we used a set of 112 reference human and NHP  
158 HBV sequences (Dataset 2, see Methods). An ML phylogenetic tree based on these  
159 reference sequences and all ancient sequences was constructed (Extended Data Fig.  
160 4). Regression of root-to-tip genetic distances against sampling dates, as well as date  
161 randomisation tests, showed a clear temporal signal in the data (Extended Data Fig. 5,  
162 SI Figs. 1-3), suggesting that molecular clock models can be applied. A dated  
163 coalescent phylogeny was constructed using BEAST2<sup>26</sup> (Fig. 2). The molecular clock  
164 was calibrated using tip dates. Strict and relaxed lognormal molecular clocks were  
165 tested with coalescent constant, exponential, and Bayesian skyline population priors  
166 (Extended Data Table 3a). Model comparisons favoured a relaxed molecular clock  
167 model with lognormally distributed rate variation and a coalescent exponential  
168 population prior (Extended Data Table 3a). The median root age of the resulting tree  
169 is estimated to 11.6 kya (95% Highest Posterior Density (HPD) interval: 8.6 to 15.3  
170 kya) and the median clock rate is  $1.18 \times 10^{-5}$  s/s/y (95% HPD interval:  $9.21 \times 10^{-6}$  to  
171  $1.45 \times 10^{-5}$  s/s/y). Under a strict molecular clock, a coalescent Bayesian skyline  
172 population prior was favoured, in which case the median root age is 15.6 kya (95%  
173 HPD interval: 13.7 to 17.8 kya) and the median substitution rate  $9.48 \times 10^{-6}$  s/s/y (95%  
174 HPD interval:  $8.3 \times 10^{-6}$  to  $1.07 \times 10^{-5}$  s/s/y) (Extended Data Tables 3a-c).

175

176 Under all model parameterisations used here, the substitution rate we find is lower  
177 than rates estimated from phylogenies built using either modern heterochronous  
178 sequences<sup>11</sup> or sequences from mother-to-child transmissions<sup>27</sup>, but higher than rates



179 inferred using external calibrations based on human migrations<sup>12</sup>. A lower rate is  
180 consistent with Tedder et al. (2013)<sup>28</sup>, who found that although mutation rates may be  
181 high, mutations within an individual often revert back to the genotype consensus, and  
182 thus rarely lead to long-term sequence change. It is also consistent with the so-called  
183 time-dependent rate phenomenon, observed for many viruses, which shows that short-  
184 term evolutionary rates are higher than long-term rates<sup>29</sup>.

185

186 The knowledge of ancient HBV genomes enables us to formally evaluate hypotheses  
187 concerning HBV origins using path sampling of calibrated phylogenies based on  
188 appropriate external divergence date assumptions. We tested several calibration points  
189 implied by a co-expansion of HBV with humans after leaving Africa for support of  
190 congruence between migrations and geographical locations of HBV clades<sup>12</sup>. We find  
191 weak evidence for a split of the F/H clade between 13.4 and 25.0 kya under a strict,  
192 but not a relaxed clock model. We do not find support for the divergence of  
193 subgenotype C3 strains between 5.1-12.0 kya, leading to a distribution in different  
194 regions of Polynesia, or for divergence of Haitian A3 strains from other genotype A  
195 strains between 0.2-0.5 kya under either strict or relaxed clock models (Extended Data  
196 Table 3d).

197

198 In the dated coalescent phylogeny, four ancient sequences (from youngest to oldest:  
199 HBV-DA119, -DA195, -RISE386, and -RISE387) group with genotype A. The first  
200 three fall well within the 7.5% nucleotide divergence criterion used to delimit  
201 membership in HBV genotypes. HBV-RISE387 is right on this limit (Extended Data  
202 Table 4a)<sup>14</sup>. The three oldest samples lack a six nucleotide insertion at the carboxyl  
203 end of the Core gene that is present in all modern genotype A viruses (Table 2)<sup>9</sup>.

204 HBV-RISE387 encodes a stop codon in its pre-Core peptide that would have ablated  
205 the expression of the immune modulator HBe antigen (HBeAg), a phenomenon  
206 known in modern HBV infections (Table 2). This characteristic viral mutant is usually  
207 found in chronic HBV carriers who seroconverted from HBeAg to anti-HBe.  
208 Interestingly, RISE386 and RISE387 are archaeologically dated only ~100 years apart  
209 and both come from the Bulanovo site in Russia, but their viruses show only 93.34%  
210 sequence identity (Extended Data Table 4b), indicating the existence of significant  
211 localized HBV diversity ~4.2 kya.

212

213 The ancient sequence HBV-DA45 phylogenetically groups with genotype B and has  
214 97.65% sequence identity with modern genotype B (Extended Data Table 4a).

215

216 Sequences HBV-DA51, -DA27, -DA222, and -DA29 phylogenetically group with the  
217 modern genotype D. They have high sequence identity (96.99 to 98.74%) with  
218 modern genotype D sequences (Extended Data Table 4a), and have the typical 33  
219 nucleotide deletion in the PreS1 sequence of the S-gene, encoding the three HBV  
220 surface proteins<sup>9</sup> (Table 2).

221

222 Sequences HBV-RISE154, -RISE254, and -RISE563 are in sister relationship to the  
223 Chimpanzee/Gorilla HBV clade (Fig. 2). HBV-RISE254 and -RISE563 have the same  
224 33 nucleotide deletion in the PreS1 sequence that is shared with NHP HBV and  
225 human genotype D (Table 2). HBV-RISE563 does not encode a functional pre-Core  
226 peptide (Table 2). Based on sequence similarity across the whole genome, HBV-  
227 RISE563 and -RISE254 together might be classified as a new human genotype that is  
228 extinct today, and HBV-RISE154 as possibly another (Extended Data Table 4).

229 However, HBV-RISE154 has low genome coverage, which precludes an exact  
230 calculation. The sister relationship of these three sequences with modern Chimpanzee  
231 and Gorilla HBV could be interpreted as a consequence of relatively recent  
232 transmission(s) of HBV from humans to NHPs<sup>15</sup>. However, other scenarios and  
233 confounding factors are possible, as these lineages are deeply separated in the tree.  
234 Incomplete lineage sorting combined with viral extinction (possibly boosted by  
235 massive recent reductions in great ape populations) should be considered. More data  
236 on current and, if possible, ancient NHP-associated HBV will be necessary to reach  
237 definitive conclusions.

238

239 The geographic locations of some of the ancient virus genotypes do not match the  
240 present-day genotype distribution, and also do not match dates and/or locations  
241 inferred in previous studies of HBV. While it is important to keep in mind that the  
242 data presented here are limited, they provide important spatiotemporal reference  
243 points in the evolutionary history of HBV. Their synopsis suggests a more  
244 complicated ancestry of present-day genotypes than previously assumed, especially in  
245 light of recent insights into the history of human migration.

246

247 We find genotype A in South-Western Russia by 4.3 kya (RISE386, RISE387), in  
248 individuals belonging to the Sintashta culture and in a sample (DA195) from the  
249 Scythian culture. The western Scythians are related to the Bronze Age cultures of the  
250 Western Steppe populations<sup>2</sup> and their shared ancestry suggests that the modern  
251 genotype A may descend from this ancient Eurasian diversity and not, as previously  
252 hypothesized, from African ancestors<sup>30,31</sup>. This is also consistent with the phylogeny  
253 (Fig. 2), as well as the fact that the three oldest ancient genotype A sequences (HBV-

254 DA195, -RISE386, and -RISE387) lack the six nucleotide insertion found in the  
255 youngest (HBV-DA119), and all modern genotype A sequences. The ancestors of  
256 subgenotypes A1 and A3 could have been carried into Africa subsequently, via  
257 migration from western Eurasia<sup>32</sup>.

258

259 The ancient HBV genotype D sequences were all found in Central Asia. HBV-DA27,  
260 found in Kazakhstan and dated to 1.6 kya, falls basal to the modern subgenotype D5  
261 sequences that today are found in the Paharia tribe from eastern India<sup>33</sup>. DA27 and the  
262 Paharia people in India are linked by their Tibeto-Burman ancestry<sup>2,34</sup>, possibly  
263 explaining the similar viruses.

264

265 Based on the observation that genotypes go extinct and can be created by  
266 recombination, our data show that the diversity we observe today is only a subset of  
267 the diversity that has ever existed. Our data support a scenario in which all present  
268 day HBV diversity arose only after the split of the Old World and New World  
269 genotypes (25-13.4 kya). Any attempt to interpret the currently known HBV tree  
270 based on human migrations that happened before this event will necessarily result in  
271 anomalies that cannot be reconciled, such as the basal position of genotypes F/H and  
272 the apical position of subgenotype C4, which is exclusively found in indigenous  
273 Australians<sup>9</sup>. If HBV did co-evolve with ancient modern humans as they left Africa as  
274 proposed previously<sup>7</sup>, most of the pattern of earlier diversity has been replaced by  
275 changes that happened after the split of the Old and New World genotypes. Genotypes  
276 F and H would therefore be remnants of the earlier now-extinct diversity, and the  
277 arrival of subgenotype C4 in Australia would have taken place long after the Old/New  
278 World split, as supported by our tree in Figure 2. Alternatively, there could have been

279 a New World origin of HBV or the introduction of HBV into humans from a different  
280 host. Our data do not allow us to speculate either way.

281

282 To our knowledge, we report the oldest exogenous viral sequences recovered from  
283 DNA of humans or any vertebrate. We show for the first time that is possible to  
284 recover viral sequences from samples of this age. We show that humans throughout  
285 Eurasia were widely infected with HBV for thousands of years. Despite the age of the  
286 samples and the imperfect diagnostic test, our dataset contained a surprisingly high  
287 proportion of HBV-positive individuals. The actual ancient prevalence during the  
288 Bronze Age and thereafter might have been higher, reaching or exceeding the  
289 prevalence typically found in contemporary indigenous populations<sup>6</sup>. This clearly  
290 establishes the potential of HBV as powerful proxy tool for research into human  
291 spread and interactions. The ancient data reveal aspects of complexity in HBV  
292 evolution that are not apparent when only modern sequences are considered. They  
293 show the existence of ancient HBV genotypes in locations incongruent with their  
294 present-day distribution, contradicting previously suggested geographic or temporal  
295 origins of genotypes or sub-genotypes; evidence for the creation of genotype A via  
296 recombination and the emergence of the genotype outside Africa; at least one now-  
297 extinct human genotype; ancient genotype-level localized diversity; and demonstrate  
298 that the viral substitution rate obtained from modern heterochronously sampled  
299 sequences is misleading. These suggest that the difficulty in formulating a coherent  
300 theory for the origin and spread of HBV may be due to genetic evidence of an earlier  
301 evolutionary scenario being overwritten by relatively recent alterations, as also  
302 suggested by Simmonds et al., in the context of recombination<sup>25</sup>. The lack of ancient  
303 sequences limits our understanding of the evolution of HBV and, very likely, of other

304 viruses. Discovery of additional ancient viral sequences may provide a clearer picture  
305 of the true origin and early diversification of HBV, enable us to address questions of  
306 paleo-epidemiology, and broaden our understanding of the contributions of natural  
307 and cultural changes (including migrations and medical practices) to human disease  
308 burden and mortality.

309

310

311 References

312

- 313 1 Allentoft, M. E. *et al.* Population genomics of Bronze Age Eurasia. *Nature*  
 314 **522**, 167-172, doi:10.1038/nature14507 (2015).
- 315 2 de Barros Damgaard, P. & Willerslev, E. 137 ancient human genomes from  
 316 across the Eurasian steppe. *Nature* (2018, in principle accepted).
- 317 3 Lai, C. L., Ratziu, V., Yuen, M.-F. & Poynard, T. Viral hepatitis B. *Lancet* **362**,  
 318 2089-2094, doi:10.1016/S0140-6736(03)15108-2 (2003).
- 319 4 Schweitzer, A., Horn, J., Mikolajczyk, R. T., Krause, G. & Ott, J. J.  
 320 Estimations of worldwide prevalence of chronic hepatitis B virus infection: a  
 321 systematic review of data published between 1965 and 2013. *Lancet* **386**,  
 322 1546-1555, doi:10.1016/S0140-6736(15)61412-X (2015).
- 323 5 *World Health Organization, Hepatitis B fact sheet*,  
 324 <<http://www.who.int/mediacentre/factsheets/fs204/en/>> (2017).
- 325 6 Murhekar, M. V., Murhekar, K. M. & Sehgal, S. C. Epidemiology of hepatitis B  
 326 virus infection among the tribes of Andaman and Nicobar Islands, India.  
 327 *Trans. R. Soc. Trop. Med. Hyg.* **102**, 729-734,  
 328 doi:10.1016/j.trstmh.2008.04.044 (2008).
- 329 7 Locarnini, S., Littlejohn, M., Aziz, M. N. & Yuen, L. Possible origins and  
 330 evolution of the hepatitis B virus (HBV). *Semin. Cancer Biol.* **23**, 561-575,  
 331 doi:10.1016/j.semcancer.2013.08.006 (2013).
- 332 8 Littlejohn, M., Locarnini, S. & Yuen, L. Origins and Evolution of Hepatitis B  
 333 Virus and Hepatitis D Virus. *Cold Spring Harb. Perspect. Med.* **6**, a021360,  
 334 doi:10.1101/cshperspect.a021360 (2016).
- 335 9 Kramvis, A. Genotypes and genetic variability of hepatitis B virus.  
 336 *Intervirolgy* **57**, 141-150, doi:10.1159/000360947 (2014).
- 337 10 Hannoun, C., Horal, P. & Lindh, M. Long-term mutation rates in the hepatitis  
 338 B virus genome. *J. Gen. Virol.* **81**, 75-83, doi:10.1099/0022-1317-81-1-75  
 339 (2000).
- 340 11 Zhou, Y. & Holmes, E. C. Bayesian estimates of the evolutionary rate and age  
 341 of hepatitis B virus. *J. Mol. Evol.* **65**, 197-205, doi:10.1007/s00239-007-0054-  
 342 1 (2007).
- 343 12 Paraskevis, D. *et al.* Dating the origin of hepatitis B virus reveals higher  
 344 substitution rate and adaptation on the branch leading to F/H genotypes. *Mol.*  
 345 *Phylogenet. Evol.* **93**, 44-54, doi:10.1016/j.ympev.2015.07.010 (2015).
- 346 13 Zehender, G. Enigmatic origin of hepatitis B virus: An ancient travelling  
 347 companion or a recent encounter? *World J. Gastroenterol.* **20**, 7622,  
 348 doi:10.3748/wjg.v20.i24.7622 (2014).
- 349 14 Kramvis, A. *et al.* Relationship of serological subtype, basic core promoter  
 350 and precore mutations to genotypes/subgenotypes of hepatitis B virus. *J.*  
 351 *Med. Virol.* **80**, 27-46, doi:10.1002/jmv.21049 (2008).
- 352 15 MacDonald, D. M., Holmes, E. C., Lewis, J. C. & Simmonds, P. Detection of  
 353 hepatitis B virus infection in wild-born chimpanzees (*Pan troglodytes verus*):  
 354 phylogenetic relationships with human and other primate genotypes. *J. Virol.*  
 355 **74**, 4253-4257 (2000).
- 356 16 Nielsen, R. *et al.* Tracing the peopling of the world through genomics. *Nature*  
 357 **541**, 302-310, doi:10.1038/nature21347 (2017).
- 358 17 Rasmussen, S. *et al.* Early divergent strains of *Yersinia pestis* in Eurasia  
 359 5,000 years ago. *Cell* **163**, 571-582, doi:10.1016/j.cell.2015.10.009 (2015).

- 360 18 Feldman, M. *et al.* A High-Coverage *Yersinia pestis* Genome from a Sixth-  
361 Century Justinianic Plague Victim. *Mol. Biol. Evol.* **33**, 2911-2923,  
362 doi:10.1093/molbev/msw170 (2016).
- 363 19 Reid, A. H., Fanning, T. G., Hultin, J. V. & Taubenberger, J. K. Origin and  
364 evolution of the 1918 "Spanish" influenza virus hemagglutinin gene. *Proc.*  
365 *Natl. Acad. Sci. U. S. A.* **96**, 1651-1656 (1999).
- 366 20 Duggan, A. T. *et al.* 17(th) Century Variola Virus Reveals the Recent History  
367 of Smallpox. *Curr. Biol.* **26**, 3407-3412, doi:10.1016/j.cub.2016.10.061 (2016).
- 368 21 Kahila Bar-Gal, G. *et al.* Tracing hepatitis B virus to the 16th century in a  
369 Korean mummy. *Hepatology* **56**, 1671-1680, doi:10.1002/hep.25852 (2012).
- 370 22 Patterson Ross, Z. *et al.* The paradox of HBV evolution as revealed from a  
371 16th century mummy. *PLOS Pathogens* **14**, e1006750,  
372 doi:10.1371/journal.ppat.1006750 (2018).
- 373 23 Bond, W. W. *et al.* Survival of hepatitis B virus after drying and storage for  
374 one week. *Lancet* **1**, 550-551 (1981).
- 375 24 Rasmussen, M. *et al.* Ancient human genome sequence of an extinct Palaeo-  
376 Eskimo. *Nature* **463**, 757-762, doi:10.1038/nature08835 (2010).
- 377 25 Simmonds, P. & Midgley, S. Recombination in the genesis and evolution of  
378 hepatitis B virus genotypes. *J. Virol.* **79**, 15467-15476,  
379 doi:10.1128/JVI.79.24.15467-15476.2005 (2005).
- 380 26 Bouckaert, R. *et al.* BEAST 2: a software platform for Bayesian evolutionary  
381 analysis. *PLoS Comput. Biol.* **10**, e1003537,  
382 doi:10.1371/journal.pcbi.1003537 (2014).
- 383 27 Simmonds, P. Reconstructing the origins of human hepatitis viruses. *Philos.*  
384 *Trans. R. Soc. Lond. B Biol. Sci.* **356**, 1013-1026, doi:10.1098/rstb.2001.0890  
385 (2001).
- 386 28 Tedder, R. S., Bissett, S. L., Myers, R. & Ijaz, S. The 'Red Queen' dilemma –  
387 running to stay in the same place: reflections on the evolutionary vector of  
388 HBV in humans. *Antivir. Ther.* **18**, 489-496, doi:10.3851/imp2655 (2013).
- 389 29 Duchêne, S., Holmes, E. C. & Ho, S. Y. W. Analyses of evolutionary  
390 dynamics in viruses are hindered by a time-dependent bias in rate estimates.  
391 *Proceedings of the Royal Society B: Biological Sciences* **281** (2014).
- 392 30 Zehender, G. *et al.* Reliable timescale inference of HBV genotype A origin  
393 and phylodynamics. *Infect. Genet. Evol.* **32**, 361-369,  
394 doi:10.1016/j.meegid.2015.03.009 (2015).
- 395 31 Hannoun, C. Phylogeny of African complete genomes reveals a West African  
396 genotype A subtype of hepatitis B virus and relatedness between Somali and  
397 Asian A1 sequences. *J. Gen. Virol.* **86**, 2163-2167, doi:10.1099/vir.0.80972-0  
398 (2005).
- 399 32 Pickrell, J. K. *et al.* Ancient west Eurasian ancestry in southern and eastern  
400 Africa. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 2632-2637,  
401 doi:10.1073/pnas.1313787111 (2014).
- 402 33 Ghosh, S. *et al.* Unique Hepatitis B Virus Subgenotype in a Primitive Tribal  
403 Community in Eastern India. *J. Clin. Microbiol.* **48**, 4063-4071,  
404 doi:10.1128/jcm.01174-10 (2010).
- 405 34 Basu, A., Sarkar-Roy, N. & Majumder, P. P. Genomic reconstruction of the  
406 history of extant populations of India reveals five distinct ancestral  
407 components and a complex structure. *Proceedings of the National Academy*  
408 *of Sciences* **113**, 1594-1599, doi:10.1073/pnas.1513197113 (2016).



410

411 **Supplementary Information** is linked to the online version of the paper at

412 [www.nature.com/nature](http://www.nature.com/nature).

413

414 **Acknowledgements** BB dedicates this work to his late mother, D. Tserendulam. We thank  
415 Stuart Rankin and the staff of the University of Cambridge High Performance Computing  
416 service and the National High-throughput Sequencing Centre (Copenhagen). This work was  
417 supported by: The Danish National Research Foundation, The Danish National Advanced  
418 Technology Foundation (The Genome Denmark platform, grant 019-2011-2), The Villum  
419 Kann Rasmussen Foundation, KU2016, European Union FP7 programme ANTIGONE (grant  
420 agreement No. 278976), European Union Horizon 2020 research and innovation  
421 programmes, COMPARE (grant agreement No. 643476), VIROGENESIS (grant agreement  
422 No. 634650). The National Reference Center for Hepatitis B and D Viruses is supported by  
423 the German Ministry of Health via the Robert Koch Institute, Berlin, Germany. BB was  
424 supported by Taylor Famil-Asia Foundation Endowed Chair in Ecology and Conservation  
425 Biology.

426

#### 427 **Author Contributions**

428 All authors contributed to the interpretation of the results.

429 BM, TJ, PD, MA, SR, MS, LO, LV, DS, DG, RF, CD, EW wrote the paper.

430 BM, TJ: screened and analysed data, created display items.

431 PD, MA: conducted sampling and generated sequence data.

432 IS, AL, EU, IP, BB, TB, KT, VM, NL, DV, EK, AE, DP, MV, TDP, VM, VS: excavated,  
433 curated, and analysed samples and archaeological context.

434 AH: designed virus capture probes.

435 LV: designed virus capture probes, performed TaqMan PCR and target enrichment  
436 experiments.  
437 AO: initiated and provided critical input on the development of NGS bioinformatics tools.  
438 DS, DG, RF: computational analysis.  
439 CD: analysed data, PCR probe design.  
440 KS, KK: conducted sampling and archaeological background.  
441 EW: initiated the work, led sampling and generation of the sequence data.

442

#### 443 **Author Information**

444 Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). The authors  
445 declare no competing financial interests. Correspondence and requests for materials should be  
446 addressed to E. W. ([ewillerslev@snm.ku.dk](mailto:ewillerslev@snm.ku.dk)).

447

448

449 **Tables**

450

451 **Table 1: Overview of samples used for phylogenetic analyses**

Sample	<sup>14</sup> C age (standard deviation)	Median cal BP age, or estimate (years)	Approx. sample age (years)	Site	Culture or period	Sex	Reads included in consensus	Coverage Consensus	Depth
RISE563	3955 (35)	4421	4488	Osterhofen- Altenmarkt, Germany	Bell Beaker	M	4383	100%	79.3x
DA222	N/D	1200-1000	1167	Butakty, Kazakhstan	Karluk	M	4132	100%	89.2x
DA195	2479 (35)	2578	2645	Sandorfalva-Eperjes, Hungary	Hungarian Scythian	F	1445	99.9%	29.2x
DA51	2220 (37)	2230	2297	Keden, Kyrgyzstan	Saka	M	712	99.2%	14.5x
RISE254	3631 (29)	3942	4009	Százhalombatta- Földvár, Hungary	Vatya	M	1491	99.0%	36.6x
DA119	N/D	1500	1567	Poprad, Slovakia	North Carpathian	M	2597	98.8%	53.1x
RISE386	3758 (34)	4121	4188	Bulanovo, Russia	Sintashta	M	331	97.8%	7.0x
DA27	1641 (33)	1543	1610	Halvay 3, Kazakhstan	Hun- Sarmatian	M	890	90.0%	14.3x
DA29	849 (25)	755	822	Karasyur, Kazakhstan	Medieval	M	222	87.5%	4.8x
DA45	2083 (27)	2053	2120	Omnogobi, Mongolia	Xiongnu	M	215	87.2%	4.3x
RISE387	3822 (33)	4215	4282	Bulanovo, Russia	Sintashta	N/ D	284	86.6%	6.2x
RISE154	3522 (24)	3784	3851	Szczepankowice, Poland	Unetice	F	128	57.2%	2.0x

452 Samples included in phylogenetic analysis, by decreasing genome coverage. Criteria for inclusion were at least  
453 50% genome coverage and sufficient aDNA damage patterns after capture. The read count indicates the number  
454 of reads used to make consensus sequences. N/D (not determined) indicates samples where dating was not  
455 performed or where osteological sex was undetermined. See Methods for information on sequence matching,  
456 consensus making, and sample dating.

457

458 **Table 2: Genome properties of ancient sequences included in phylogenetic analyses**

Sample	Genotype of closest sequence	Sequence identity to closest sequence	Genome length	Sero-type	Insertions / deletions	Predicted HBeAg status
DA119	A3	97.8%	3221	<i>adw2</i>	6nt insert at the C-terminus of core region	Positive
DA195	A3	96.2%	3215	<i>adw2</i>	None	Positive
RISE386	A	95.2%	3215	<i>adw2</i>	None	Positive
RISE387	A	92.5%	3215	<i>adw2</i>	None	Negative PreC stop codon
DA45	B1	96.6%	3215	<i>ayw1</i>	None	Positive
DA29	D3	98.5%	3182	<i>ayw2</i>	33nt deletion at the N-terminus of the preS1 region	Positive
DA222	D3	98.7%	3182	<i>ayw2</i>	33nt deletion at the N-terminus of the preS1 region	Positive
DA27	D1	97.2%	3182	<i>ayw2</i>	33nt deletion at the N-terminus of the preS1 region	Positive
DA51	D1	96.7%	3182	<i>ayw2</i>	33nt deletion at the N-terminus of the preS1 region	Positive
RISE154	Chimp.	92.5%	Ambiguous	<i>adw2</i> *	Ambiguous	Positive
RISE254	Chimp.	95.2%	3182	<i>adw2</i>	33nt deletion at the N-terminus of the preS1 region	Positive
RISE563	Gorilla	92.7%	3182	<i>adw2</i>	33nt deletion at the N-terminus of the preS1 region	Negative PreC stop codon

459 Genotype groups are sorted by increasing sample age. \* Serotype could not be determined unambiguously, due  
 460 to lack of coverage.

461

## 462 **Figure legends**

### 463 **Figure 1: Geographic distribution of analysed samples and modern genotypes**

464 **a**, Distribution of modern human HBV genotypes<sup>8</sup>. Genotypes relevant to the manuscript are shown in colour.  
465 Coloured shapes indicate the locations of the HBV-positive samples included for further analysis, as in panel **b**.  
466 **b**, Locations of analysed Bronze Age samples<sup>1</sup> are shown as circles, Iron Age and later samples<sup>2</sup>, as triangles.  
467 Coloured markers indicate HBV-positive samples. Ancient genotype A samples are found in regions where  
468 genotype D predominates today, and DA27 is of sub-genotype D5 which today is found almost exclusively in  
469 India.

470

### 471 **Figure 2: Dated maximum clade credibility tree of HBV**

472 A lognormal relaxed clock and coalescent exponential population prior were used. Grey horizontal bars indicate  
473 the 95% HPD interval of the age of the node. Larger numbers on the nodes indicate the age and 95% HPD  
474 interval of the age under a strict clock and Bayesian skyline tree prior. Clades of genotypes C (except clade C4),  
475 F, and H are collapsed and shown as dots. Taxon names indicate: genotype / subgenotype, accession number,  
476 sample age, country abbreviation of sequence origin, region of sequence origin, host species, and optional  
477 additional remarks.

478

## 479 **Methods**

### 480 **HBV datasets**

481 The following HBV datasets were used in the present study. Full listings of accession  
482 numbers are given in the Supplementary Methods.

483 **Dataset 1:** 26 HBV genomes, covering all species in the *Orthohepadnaviridae*. This includes  
484 one sequence each from the human HBV genotypes (A-J), Orangutan, Chimpanzee, Gorilla,  
485 Gibbon, Woolly monkey, Woodchuck, Ground squirrel, Arctic ground squirrel, Horseshoe  
486 bat, and four sequences from Roundleaf bats and three from Tent-making bats, largely  
487 following Drexler et al.<sup>35</sup>

488 **Dataset 2:** 124 HBV genomes, from humans and NHP. This set contains 92 sequences from  
489 Paraskevis et al.<sup>12</sup> (excluding their incomplete sequences), 7 additional genotype D  
490 sequences, the Korean mummy genotype C sequence<sup>21</sup>, the 12 ancient sequences from the  
491 present study, and 12 full genomes selected from a set of 9066 full HBV genomes  
492 downloaded from NCBI<sup>36</sup> on 2017-08-24 (Entrez query: hepatitis b virus[organism] not  
493 rna[title] not clone[title] not clonal[title] not patent[title] not recombinant[title] not  
494 recombination[title] and 3000:4000[sequence length]) corresponding to the closest, non-  
495 artificial match for each of the ancient sequences. Dates for these sequences were acquired by  
496 looking for a date of sample collection in the NCBI entry, or the paper where the sequence  
497 was first published. If a range of dates was mentioned, the mean was used. If no date of  
498 sample collection was found in this way, either the year of the publication of the paper, or the  
499 year of addition of the sequence to GenBank was used, whichever was earlier.

500 **Dataset 3:** 124 HBV genomes, from humans, NHP, and a variety of other Orthohepadnavirus  
501 host species, including Woolly monkey, Roundleaf and Tent-making bat, Ground and Arctic  
502 ground squirrel, Woodchuck, and Snow goose. This set contains 113 sequences that are the

503 union of a selection of 91 sequences from Paraskevis et al.<sup>12</sup> and 29 from Drexler et al.<sup>35</sup>, and  
504 11 additional sequences.

505 **Dataset 4:** 3505 HBV genomes. 3384 are from Bell et al., (2016)<sup>37</sup>, divided into ten human  
506 genotypes. To these we added 17 Chimpanzee, 56 Gorilla, 12 Gibbon and 36 Orangutan full  
507 HBV genome sequences downloaded from NCBI on 2017-01-18, resulting in 14 genome  
508 categories.

### 509 **Dating of ancient samples**

510 Sample ages were determined by direct <sup>14</sup>C-dating. These ages were calibrated using OxCal<sup>38</sup>  
511 (version 4.3) using the IntCal13 curve<sup>39</sup>. Table 1 shows the <sup>14</sup>C age and standard deviation for  
512 each sample. This is followed by the median probability calibrated age before present (cal  
513 BP), where “present” is defined as 1950. RISE386 was <sup>14</sup>C dated twice, with ages (standard  
514 deviation) of 3740 (33) and 3775 (34), so a rounded mean of 3758 (34) was used for its  
515 calibration. DA29 was dated at 822 years using <sup>14</sup>C and also at ~700 years using multi-proxy  
516 methods, the former was used for consistency. The dates for DA119, DA222, RISE548,  
517 RISE556, RISE568, and RISE597 are best estimates, based on sample context.

### 518 **Data and data processing**

519 We analysed 101 Bronze Age samples published in Allentoft, et al.<sup>1</sup>, 137 predominantly Iron  
520 Age samples published in Damgaard et al.<sup>2</sup>, and 66 additional samples from the Bronze Age.  
521 A total of 114.58x10<sup>9</sup> Illumina HiSeq 2500 sequencing reads were processed.

522

523 AdapterRemoval<sup>40</sup> (version 2.1.7) was used with its default settings to remove adaptors from  
524 all sequences, to trim N bases from the ends of reads, and to trim bases with quality ≤ 2.

525 Reads were aligned against a human genome (GRCh38<sup>41</sup>) using BWA<sup>42</sup> (version 0.7.15-



526 r1140, mem algorithm). Reads that did not match the human genome were then mapped  
527 against the NCBI viral protein reference database containing 274,038 viral protein sequences  
528 (downloaded on 2016-08-31) using DIAMOND<sup>43</sup> (version 0.8.25). Protein matches were  
529 grouped into their corresponding viruses. Reads matching HBV were found in 25 samples.  
530

531 The non-human reads from the HBV-positive samples that had more than three reads  
532 matching HBV using DIAMOND were selected for a subsequent BLAST<sup>44</sup> (version 2.4.0)  
533 analysis. A BLAST database was made from Dataset 3, and samples were matched using  
534 blastn (with arguments -task blastn -evalue 0.01). Matching reads with bit scores greater than  
535 50 for all samples (except DA222 (70) and DA45 (55)) were selected for subsequent  
536 processing. The number of reads selected from the BLAST matches, per sample, is shown in  
537 Table 1, with additional detail in Extended Data Table 1. Across all samples 11,149 reads  
538 matched against HBV sequences.

### 539 **PCR confirmation**

540 Real-time PCR was established using primers and TaqMan probes as described by Drosten et  
541 al.,<sup>45</sup> which amplifies a 91 base pair amplicon of the HBV genome. Primers and probe were  
542 added to QuantiTect PCR mix (Qiagen #204343) in a final concentration of 400 nM or 200  
543 nM, respectively, in a total reaction volume of 25 ul, including 5 ul template. Using the  
544 Roche LC480 or Agilent Mx3006p instruments, PCRs were incubated for 15 min. at 95°C  
545 followed by 45 cycles of 15 seconds at 94°C and 60 seconds at 60°C, measuring fluorescence  
546 from the 6-carboxy-fluorescein/BHQ1-labelled probe and the passive dye (ROX) at the end  
547 of each cycle.

548 Careful precautions were taken to prevent PCR contamination. PCR mastermixes were  
549 prepared in dedicated ancient DNA clean lab facilities, in which no prior targeted work has  
550 been carried out on HBV. Ancient DNA extracts and non-template controls (NTC) were

551 added into PCR reactions in this location too, which were not subsequently opened. Positive  
552 control material was handled in labs in a physically separated building. Here, standard  
553 material, diluted to 5-50 copies/reaction, was added to duplicate PCR reactions along with  
554 additional NTCs.

## 555 **Virus capture**

556 14 samples with sufficient sample material were selected for virus capture (DA27, DA29,  
557 DA45, DA51, DA85, DA89, DA119, DA195, DA222, RISE254, RISE386, RISE416,  
558 RISE568, RISE556). The viral reference genomes for probes were selected as follows. The  
559 International Committee for Taxonomy of Viruses (ICTV) 2012 listed 2618 viral species. As  
560 many had no associated reference genomes or merely partial sequence information, we  
561 selected 2599 sequences of full-length viral genomes, available from GenBank (June 2014),  
562 representing viral species found in vertebrates excluding fish. Sequences <1000 nt were  
563 discarded. Sequences with identical length and organism ID were regarded as duplicates and  
564 thus reduced to 1. For a number of specific viral taxa for which a large number of similar  
565 reference sequences are available, we manually selected representative genomes or genome  
566 segments (SI Tables 1 and 2). For example, among 72 available Hepatitis C virus genome  
567 sequences, we selected one genome per subtype (1a-c, g; 2a-c, i, k; 3a, b, i, k; 4a-d, f, g, k-r,  
568 t; 5a; 6a-u; 7a). Likewise, 12 HIV-1 genomes were selected representing groups M (subtypes  
569 A-D, F1, F2, H, J, K, N, O, and P). For influenza A virus, we included only sequences from  
570 segment 7 and segment 5 encoding the conserved matrix proteins M1/M2 and the  
571 nucleocapsid protein NP, respectively. We selected 82 M1/M2 segments and 115 NP  
572 segments among the available segments sequences. All available segments were included  
573 from genomes belonging to *Arenaviridae*, *Bunyaviridae*, and *Reoviridae*. For members of  
574 *Poxvirinae* for which full genomes were unavailable (Skunk-, Raccoon-, and Volepox virus)  
575 sequences representing the conserved gene encoding the DNA-dependent RNA polymerase

576 were included (n=22). In addition, 2 partial genomes of Squirrelpox virus were included. By  
577 mistake 2 and 9 partial sequences were included from *Iridoviridae* (1.5-2.5 kb) and  
578 *Coronaviridae* (1.3-14.5 kb), respectively, already represented by full genomes. Likewise,  
579 sequences representing Merkel cell polyomavirus and KI polyomavirus were not included  
580 among the reference genomes used for probe design. SeqCap EZ hybridization probes were  
581 designed and synthesized by Roche NimbleGen (Madison, USA) based on the resulting  
582 reference sequences.

583 Capture was performed on double-indexed libraries prepared from ancient DNA, following  
584 the manufacturer's protocol (version 4.3) with the following modifications. Briefly, 1.8 to 2.2  
585 µg of pooled libraries were hybridized at 47°C for 65-70 hours with low complexity C<sub>0</sub>T-1  
586 DNA, specific P5/P7 adaptor-blocking oligonucleotides each containing a hexamer motif of  
587 inosine nucleotides to match individually indexed adapters, hybridization buffer containing  
588 10% formamide, and the capture probes. Dynabeads M-270 (Invitrogen) were used to recover  
589 the hybridized library fragments. After washing and eluting the libraries, the post-capture  
590 PCR amplification was performed with KAPA Uracil+ polymerase (Kapa Biosystems). PCR  
591 cycling conditions were as follows: 1 cycle of 3 min at 95°C, followed by 14 cycles of: 20  
592 sec denaturation at 98°C, 15 sec annealing at 65°C and 30 sec elongation at 72 °C, ending  
593 with 5 min at 72°C. The amplified captured libraries were purified using AMPureXP beads  
594 (Agencourt).

595 Shotgun sequencing data was generated as described in Allentoft et al. (2015)<sup>1</sup>. Sequencing  
596 of target-enriched libraries was performed on Illumina Hiseq2500 SR80bp, V4 chemistry.  
597 The resulting reads were compared to Dataset 2 using BLASTn (with arguments -task blastn -  
598 evalue 0.01). Matching reads with bit scores greater than 50 for all samples (except DA222  
599 (70) and DA45 (55)) were selected for subsequent processing. In total, 6757 reads matched  
600 HBV after capture.

## 601 **Sequence authenticity**

602 The following evidence leads us to believe that the ancient HBV sequences are authentic and  
603 that the possibility of contamination can be excluded:

- 604 (1) Standard precautions for working with ancient DNA were applied<sup>46</sup>.
- 605 (2) Sequences were checked for typical ancient DNA damage patterns using  
606 mapDamage<sup>47</sup> (version 2.0.6). Whenever sufficient amounts of data were  
607 available (>200 HBV reads), we found C>T mutations at the 5' end, typical of  
608 ancient DNA<sup>48</sup> (see Extended Data Fig. 1a,c).
- 609 (3) Capture was performed on sample DA222 DNA extracts with and without pre-  
610 treatment by Uracil-Specific Excision Reagent (USER)<sup>49</sup>. After USER treatment  
611 (3h at 37°C) of the aDNA extract, the damage pattern is eliminated (Extended  
612 Data Fig. 1b).
- 613 (3) As the ancient viruses are from three different HBV genotypes (A, B, D) and a  
614 clade in sister relationship to NHP viruses, any argument that samples were  
615 contaminated would have to account for this diversity as well as the sequence  
616 novelty.
- 617 (4) HBV sequences were identified in 25 of 305 analysed samples (Table 1), showing  
618 that the findings cannot be due to a ubiquitous laboratory contaminant.
- 619 (5) Despite the low frequency of positive samples, we sequenced extraction blanks to  
620 provide additional evidence against the possibility that the HBV sequences  
621 stemmed from sporadic incorporation, amplification, and sequencing of  
622 background reagent contaminants into the ancient DNA libraries. The negative  
623 extraction controls were amplified for 40 PCR cycles, and BLAST was used to  
624 match the read sequences against Dataset 3, with the same parameters used for the  
625 ancient samples. Because the ancient HBV positive reads used to assemble

626 genomes all had bit scores of at least 50 (see Data and Data Processing, above),  
627 we filtered the negative extraction control BLAST output for reads with a bit  
628 score  $\geq 45$ . No reads (out of 23 million) matched any HBV genome at that level.

629 (6) HBV is a blood-borne virus that is mainly transmitted by exposure to infectious  
630 blood and that does not occur in the environment<sup>3</sup>, making contamination during  
631 archaeological excavation extremely unlikely.

### 632 **Consensus sequences**

633 Reads from the original sequencing and from the capture were aligned to a reference genome  
634 (SI Table 3) in Geneious<sup>50</sup> (version 9) using Medium Sensitivity / Fast and Iterate up to 5  
635 times. Because aDNA damage often clusters towards read termini<sup>48</sup>, the resulting alignments  
636 were carefully curated by hand to remove non-matching termini of reads if the majority of the  
637 read showed a very good match with the reference sequence.

### 638 **Genotyping**

639 All reads used to construct the ancient HBV consensus sequences were matched against the  
640 full NCBI nucleotide (nt) database (downloaded December 28, 2016) using BLAST. 97.5%  
641 of the reads had HBV as their top match. All ancient consensus sequences were matched  
642 against the full HBV genomes of Dataset 4 with the Needleman-Wunsch algorithm<sup>51</sup>, as  
643 implemented in EMBOSS<sup>52</sup> (version 6.6.0.0). For each ancient sequence, the percent  
644 sequence identity for each modern genotype and four NHP species is listed in Extended Data  
645 Table 4a. The Needleman-Wunsch algorithm was also used to calculate the pairwise  
646 sequence similarity between all ancient sequences (Extended Data Table 4b).

### 647 **Recombination analysis**

648 The Recombination Detection Program<sup>53</sup>, version 4 (RDP4) was used to search for evidence  
649 of recombination within the 12 ancient sequences and a selection of 15 modern human and  
650 NHP sequences (Supplementary Methods). Recombination with HBV-RISE387 as the  
651 recombinant and HBV-DA51 as one parent, was suggested at positions 1567-2256, by seven  
652 recombination methods (RDP<sup>54</sup>, GENECONV<sup>55</sup>, BootScan<sup>56</sup>, MaxChi<sup>57</sup>, Chimaera<sup>58</sup>,  
653 SiScan<sup>59</sup>, and 3Seq<sup>60</sup>) with p-values from  $1.179 \times 10^{-6}$  to  $5.336 \times 10^{-11}$  (Extended Data Table 2).  
654 The same recombination was suggested for all 4 ancient genotype A and two modern  
655 genotype A sequences. Graphical evidence of the recombination and the predicted break  
656 point distribution for sequences HBV-RISE386 and HBV-RISE387 from three methods  
657 (MaxChi, Bootscan, and RDP) is shown in Extended Data Fig. 3.

## 658 **Phylogenetic analysis**

### 659 **Initial maximum likelihood phylogenies**

660 An initial Maximum Likelihood (ML) tree was generated to ascertain that the ancient  
661 sequences fall within the primate HBV clades. Dataset 1 and the ancient sequences were  
662 aligned in MAFFT<sup>61</sup> (version 7). The ML tree was constructed using PhyML<sup>62</sup> (version  
663 20160116), optimizing topology, branch lengths, and rates. We used a GTR substitution  
664 model, with base frequencies determined by ML, and an ML-estimated proportion of  
665 invariant sites and 100 bootstraps (Extended Data Fig. 2). Furthermore, an ML tree (Extended  
666 Data Fig. 4) was generated based on a MAFFT alignment of Dataset 2 and the ancient  
667 sequences, using the same parameters as outlined above. The final trees show nodes with  
668 support values less than 70 as polytomies.

### 669 **Dated coalescent phylogenies**

670 In order to check for a temporal signal in the data, a root-to-tip regression and date  
671 randomisation tests were performed. For the root-to-tip regression, input trees were  
672 calculated using Dataset 2 with the addition of a Woolly Monkey sequence (GenBank

673 Accession Number: AF046996) as an outgroup. Three phylogenetic algorithms were used,  
674 Neighbour Joining, ML (PhyML), and Bayesian (MrBayes<sup>63</sup> (version 3.2.5)) methods (SI  
675 Figs. 1-3). Root-to-tip distances were extracted using TempEst<sup>64</sup> (version 1.5). For ML and  
676 Bayesian, root distances for tip taxa (in substitutions per site) were extracted from optimized  
677 tree topologies (ML and Maximum Clade Credibility trees, respectively). For NJ, root-to-tip  
678 distances were averaged over 1000 bootstrap replicates. Regression analyses were performed  
679 with Scipy<sup>65</sup> (version 0.16.0). For the date randomisation tests, we used three different  
680 approaches to randomise tip dates: First, tip dates were randomised between all sequences in  
681 the phylogeny. Second, tip dates were randomised only among the ancient sequences  
682 presented in this paper, as well as the Korean mummy sequence (accession number  
683 JN315779). The modern sequences retained their correct ages. Third, dates were randomised  
684 within a clade. For each of the three approaches, we performed three independent  
685 randomisations. This resulted in a total of nine analyses, which were run for 100,000,000  
686 generations each, under the relaxed lognormal clock model and coalescent exponential tree  
687 prior. We also ran the same analyses under a strict clock and coalescent Bayesian skyline tree  
688 prior, which were run for 20,000,000 generations. We used a GTR substitution model with  
689 unequal base frequencies, four gamma rate categories, estimated gamma distribution of rate  
690 variation, and estimated proportion of invariant sites, as found by bModelTest<sup>66</sup> (version  
691 1.0.4). None of the analyses using the relaxed clock converged (Estimated Sample Size (ESS)  
692 < 200). This is most likely because the mis-specification of the dates leads to an incongruence  
693 between the sequence and time information. Under the strict clock model, all runs converged,  
694 and none of the 95% HPD intervals of the tree height overlapped between the randomised  
695 and the non-randomised runs, fulfilling the criteria for evidence of a temporal signal<sup>67</sup>.

696 Dated phylogenies were estimated using BEAST2<sup>26</sup> (version 2.4.4, prerelease). We used a  
697 MAFFT alignment of Dataset 2. Using bModelTest<sup>66</sup>, we selected a GTR substitution model

698 with unequal base frequencies, four gamma rate categories, estimated gamma distribution of  
699 rate variation, and estimated proportion of invariant sites. Proper priors were used  
700 throughout. Path sampling, as implemented in BEAST2, was performed to select between  
701 strict or relaxed lognormal clock and a coalescent constant, exponential, or coalescent  
702 Bayesian skyline tree prior (Extended Data Table 3a). Likelihood values were compared  
703 using a Bayes factor test. According to Kass and Raftery<sup>68</sup>, a Bayes factor in the range of 3-  
704 20 implies positive support, 20-150 strong support, and >150 overwhelming support. The  
705 relaxed lognormal clock model in combination with a coalescent exponential tree prior was  
706 favoured. For the final tree, a Markov chain Monte Carlo analysis was run until parameters  
707 reached an ESS > 200, sampling every 2000 generations. Convergence and mixing were  
708 assessed using Tracer<sup>69</sup> (version 1.6). The final tree files were subsampled to contain 10,000  
709 or 10710 (for the relaxed lognormal clock, coalescent exponential tree prior) trees, with the  
710 first 25% of samples discarded as burn-in. Maximum clade credibility trees were made using  
711 TreeAnnotator<sup>26</sup> (version 2.4.4 prerelease).

712

713 In order to formally test the Out of Africa hypothesis, calibration points were tested using  
714 path sampling as implemented in BEAST2. Calibration points were constrained as follows.  
715 Split of genotypes F and H: The MRCA of all genotype F and H sequences was constrained  
716 using a uniform(13,400: 25,000) distribution, as this is the range of estimates for when the  
717 Americas were first colonized<sup>70,71</sup>. Split of subgenotype A3 in Haiti: The MRCA of FJ692598  
718 and FJ692611 was constrained using a uniform(200: 500) distribution, due to the timing of  
719 the slave trade to Haiti<sup>72</sup>. Split of C3 in Polynesia: The MRCA of X75656 and X75665 was  
720 constrained using a uniform(5,100: 12,000) distribution, due to the range of estimates for the  
721 MRCA of Polynesian populations<sup>12,73</sup>. Calibration points were tested under both a relaxed



722 lognormal clock, coalescent exponential tree prior, and a strict clock, Bayesian skyline tree  
723 prior.

724

725 **Data availability**

726 The complete sequences in this study have been deposited in the European Nucleotide

727 Archive under sample accession numbers ERS2295383-ERS2295394.

728

729

730

731 **References**

- 732 35 Drexler, J. F. *et al.* Bats carry pathogenic hepadnaviruses antigenically related to  
733 hepatitis B virus and capable of infecting human hepatocytes. *Proc. Natl. Acad. Sci.*  
734 *U. S. A.* **110**, 16151-16156, doi:10.1073/pnas.1308049110 (2013).
- 735 36 Geer, L. Y. *et al.* The NCBI BioSystems database. *Nucleic Acids Res.* **38**, D492-496,  
736 doi:10.1093/nar/gkp858 (2010).
- 737 37 Bell, T. G., Yousif, M. & Kramvis, A. Bioinformatic curation and alignment of  
738 genotyped hepatitis B virus (HBV) sequence data from the GenBank public  
739 database. *Springerplus* **5**, 1896, doi:10.1186/s40064-016-3312-0 (2016).
- 740 38 Bronk Ramsey, C. & Ramsey, C. B. Bayesian Analysis of Radiocarbon Dates.  
741 *Radiocarbon* **51**, 337-360, doi:10.1017/s0033822200033865 (2009).
- 742 39 Reimer, P. J. *et al.* IntCal13 and Marine13 Radiocarbon Age Calibration Curves 0–  
743 50,000 Years cal BP. *Radiocarbon* **55**, 1869-1887, doi:10.2458/azu\_js\_rc.55.16947  
744 (2013).
- 745 40 Lindgreen, S. AdapterRemoval: easy cleaning of next-generation sequencing reads.  
746 *BMC Res. Notes* **5**, 337, doi:10.1186/1756-0500-5-337 (2012).
- 747 41 *Human Genome Overview - Genome Reference Consortium*,  
748 <<https://www.ncbi.nlm.nih.gov/grc/human>> (2017).
- 749 42 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler  
750 transform. *Bioinformatics* **25**, 1754-1760, doi:10.1093/bioinformatics/btp324 (2009).
- 751 43 Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using  
752 DIAMOND. *Nat. Methods* **12**, 59-60, doi:10.1038/nmeth.3176 (2015).
- 753 44 Camacho, C. *et al.* BLAST : architecture and applications. *BMC Bioinformatics* **10**,  
754 421, doi:10.1186/1471-2105-10-421 (2009).
- 755 45 Drosten, C., Weber, M., Seifried, E. & Roth, W. K. Evaluation of a new PCR assay  
756 with competitive internal control sequence for blood donor screening. *Transfusion* **40**,  
757 718-724 (2000).
- 758 46 Willerslev, E. & Cooper, A. Review Paper. Ancient DNA. *Proceedings of the Royal*  
759 *Society B: Biological Sciences* **272**, 3-16, doi:10.1098/rspb.2004.2813 (2005).
- 760 47 Jónsson, H., Ginolhac, A., Schubert, M., Johnson, P. L. F. & Orlando, L.  
761 mapDamage2.0: fast approximate Bayesian estimates of ancient DNA damage  
762 parameters. *Bioinformatics* **29**, 1682-1684, doi:10.1093/bioinformatics/btt193 (2013).
- 763 48 Orlando, L., Gilbert, M. T. P. & Willerslev, E. Reconstructing ancient genomes and  
764 epigenomes. *Nat. Rev. Genet.* **16**, 395-408, doi:10.1038/nrg3935 (2015).
- 765 49 Briggs, A. W. *et al.* Removal of deaminated cytosines and detection of in vivo  
766 methylation in ancient DNA. *Nucleic Acids Res.* **38**, e87-e87,  
767 doi:10.1093/nar/gkp1163 (2010).
- 768 50 Kearse, M. *et al.* Geneious Basic: an integrated and extendable desktop software  
769 platform for the organization and analysis of sequence data. *Bioinformatics* **28**, 1647-  
770 1649, doi:10.1093/bioinformatics/bts199 (2012).
- 771 51 Needleman, S. B. & Wunsch, C. D. A general method applicable to the search for  
772 similarities in the amino acid sequence of two proteins. *J. Mol. Biol.* **48**, 443-453  
773 (1970).
- 774 52 Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open  
775 Software Suite. *Trends Genet.* **16**, 276-277 (2000).
- 776 53 Martin, D. P., Murrell, B., Golden, M., Khoosal, A. & Muhire, B. RDP4: Detection and  
777 analysis of recombination patterns in virus genomes. *Virus Evol* **1**, vev003,  
778 doi:10.1093/ve/vev003 (2015).
- 779 54 Martin, D. & Rybicki, E. RDP: detection of recombination amongst aligned  
780 sequences. *Bioinformatics* **16**, 562-563, doi:10.1093/bioinformatics/16.6.562 (2000).
- 781 55 Padidam, M., Sawyer, S. & Fauquet, C. M. Possible Emergence of New  
782 Geminiviruses by Frequent Recombination. *Virology* **265**, 218-225,  
783 doi:10.1006/viro.1999.0056 (1999).

- 784 56 Martin, D. P., Posada, D., Crandall, K. A. & Williamson, C. A Modified Bootscan  
785 Algorithm for Automated Identification of Recombinant Sequences and  
786 Recombination Breakpoints. *AIDS Res. Hum. Retroviruses* **21**, 98-102,  
787 doi:10.1089/aid.2005.21.98 (2005).
- 788 57 Smith, J. Analyzing the mosaic structure of genes. *J. Mol. Evol.* **34**,  
789 doi:10.1007/bf00182389 (1992).
- 790 58 Posada, D. & Crandall, K. A. Evaluation of methods for detecting recombination from  
791 DNA sequences: computer simulations. *Proc. Natl. Acad. Sci. U. S. A.* **98**, 13757-  
792 13762, doi:10.1073/pnas.241370698 (2001).
- 793 59 Gibbs, M. J., Armstrong, J. S. & Gibbs, A. J. Sister-Scanning: a Monte Carlo  
794 procedure for assessing signals in recombinant sequences. *Bioinformatics* **16**, 573-  
795 582, doi:10.1093/bioinformatics/16.7.573 (2000).
- 796 60 Boni, M. F., Posada, D. & Feldman, M. W. An Exact Nonparametric Method for  
797 Inferring Mosaic Structure in Sequence Triplets. *Genetics* **176**, 1035-1047,  
798 doi:10.1534/genetics.106.068874 (2006).
- 799 61 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:  
800 improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772-780,  
801 doi:10.1093/molbev/mst010 (2013).
- 802 62 Guindon, S. *et al.* New Algorithms and Methods to Estimate Maximum-Likelihood  
803 Phylogenies: Assessing the Performance of PhyML 3.0. *Syst. Biol.* **59**, 307-321,  
804 doi:10.1093/sysbio/syq010 (2010).
- 805 63 Ronquist, F. & Huelsenbeck, J. P. MrBayes 3: Bayesian phylogenetic inference  
806 under mixed models. *Bioinformatics* **19**, 1572-1574 (2003).
- 807 64 Rambaut, A., Lam, T. T., Max Carvalho, L. & Pybus, O. G. Exploring the temporal  
808 structure of heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus*  
809 *Evol* **2**, vew007, doi:10.1093/ve/vew007 (2016).
- 810 65 SciPy, <<http://www.scipy.org>> (2017).
- 811 66 Bouckaert, R. R. & Drummond, A. J. bModelTest: Bayesian phylogenetic site model  
812 averaging and model comparison. *BMC Evol. Biol.* **17**, 42, doi:10.1186/s12862-017-  
813 0890-6 (2017).
- 814 67 Duchêne, S., Duchêne, D., Holmes, E. C. & Ho, S. Y. W. The Performance of the  
815 Date-Randomization Test in Phylogenetic Analyses of Time-Structured Virus Data.  
816 *Mol. Biol. Evol.* **32**, 1895-1906, doi:10.1093/molbev/msv056 (2015).
- 817 68 Kass, R. E. & Raftery, A. E. Bayes Factors. *J. Am. Stat. Assoc.* **90**, 773,  
818 doi:10.2307/2291091 (1995).
- 819 69 Rambaut, A., Suchard, M. A., Xie, D. & Drummond, A. J. *Tracer v1.6*,  
820 <<http://beast.community/tracer>> (2017).
- 821 70 Sanchez, G. *et al.* Human (Clovis)-gomphothere (*Cuvieronius* sp.) association ~  
822 13,390 calibrated yBP in Sonora, Mexico. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 10972-  
823 10977, doi:10.1073/pnas.1404546111 (2014).
- 824 71 Bourgeon, L., Burke, A. & Higham, T. Earliest Human Presence in North America  
825 Dated to the Last Glacial Maximum: New Radiocarbon Dates from Bluefish Caves,  
826 Canada. *PLoS One* **12**, e0169486, doi:10.1371/journal.pone.0169486 (2017).
- 827 72 Andernach, I. E., Nolte, C., Pape, J. W. & Muller, C. P. Slave trade and hepatitis B  
828 virus genotypes and subgenotypes in Haiti and Africa. *Emerg. Infect. Dis.* **15**, 1222-  
829 1228, doi:10.3201/eid1508.081642 (2009).
- 830 73 Kayser, M. *et al.* Melanesian and Asian origins of Polynesians: mtDNA and Y  
831 chromosome gradients across the Pacific. *Mol. Biol. Evol.* **23**, 2234-2244,  
832 doi:10.1093/molbev/msl093 (2006).
- 833

834

835

## 836 **Extended Data table and figure titles and legends**

837

### 838 **Extended Data Table 1 | Extended overview of samples with reads matching HBV and**

#### 839 **PCR results**

840 **a**, Extended overview of samples with reads matching HBV. Rows are sorted by decreasing

841 consensus coverage. Explanation of column titles, from left to right starting from the second

842 column: <sup>14</sup>C age and standard deviation; Median cal BP age or estimate (in years);

843 approximate sample age in years; site; culture or period; gender; number of sequencing

844 reads that matched HBV using DIAMOND<sup>43</sup>; number of HBV proteins matched by those

845 reads; number of sequencing reads that matched HBV using a BLASTn<sup>44</sup> database built

846 from Dataset 3 (see Methods); the number of reads from the capture that matched HBV

847 using BLASTn (as above); the bit score cut-off above which matching reads were used to

848 form consensus sequences; the percentage of the consensus genome covered by matching

849 reads; average depth of coverage across the reference genome, as reported by Geneious<sup>50</sup>.

850 When reading sample information across a row, an empty cell will be encountered when

851 processing on that sample was concluded, either (in column 6) due to too few matching

852 reads or (penultimate column) consensus coverage less than 50%. **b**, TaqMan PCR results.

853 Four extracts from samples with HBV reads were selected for TaqMan PCR confirmation:

854 two with a large proportion of HBV reads (DA222 and DA195), two with a small proportion of

855 HBV reads (DA85 and DA89), and one with no HBV reads (DA351). HBV was detected in

856 extracts from DA222 and DA195, whereas the three low- and zero-read samples were

857 negative, as were all non-template controls.

858

### 859 **Extended Data Table 2 | Genotype A predicted recombination break points and p-**

#### 860 **values**

861 **a**, The p-values assigned to the predicted genotype A recombination by the seven methods

862 used by RDP4<sup>53</sup>, in the order given by RDP. The number of sequences in which the

863 recombination was predicted is always 6, corresponding to the 4 ancient and two modern  
864 genotype A sequences. **b**, The predicted start and end break points for each of the 6  
865 genotype A sequences. Sequences are ordered from oldest to youngest. The 99%  
866 confidence intervals for the start and end points are shown (n=15 sequences analysed in all  
867 cases), and are identical for all sequences. The predicted break points are close to the  
868 boundaries of the polymerase. For example, for the modern genotype A sequence  
869 LC074724, the polymerase is found in regions 1-1623 and 2307-3221 and the predicted  
870 break points are 1622 and 2256. If recombination formed an HBV-RISE387/6 ancestor, it is  
871 possible that the entire polymerase gene was contributed by one parent.

872

### 873 **Extended Data Table 3 | Model testing and inferred age of genotypes**

874 Models were compared using Path Sampling, as implemented in BEAST2<sup>28</sup>. Likelihood  
875 values were compared using a Bayes factor test. A positive value for the Bayes factor  
876 implies support for model 1, a negative value support for model 2. According to Kass and  
877 Raftery<sup>68</sup>, a Bayes factor in the range of 3-20 implies positive support, 20-150 strong  
878 support, and >150 overwhelming support. **a**, Results of testing different clock models and  
879 population assumptions to be used for dated phylogenies. Positive numbers indicate support  
880 for the columns model, negative number for the rows model. **b**, MRCA age of individual  
881 nodes under a strict clock and Bayesian skyline tree prior or under a relaxed lognormal clock  
882 and coalescent exponential tree prior. **c**, Root age and substitution rates under different  
883 clock models and tree priors. **d**, Results of testing different calibration point hypotheses  
884 under a strict clock and Bayesian skyline tree prior or under a relaxed lognormal clock and  
885 coalescent exponential tree prior.

886

### 887 **Extended Data Table 4 | Consensus sequence identity**

888 **a**, Best consensus sequence identity with 14 groups of HBV full genomes. The Needleman-  
889 Wunsch algorithm (as implemented in EMBOSS<sup>52</sup>) was used to globally align each sample  
890 consensus sequence against each of the 3384 full HBV genomes of Dataset 4 (see

891 Methods). The table shows the best nucleotide (nt) similarity percentage for each sample  
892 consensus against 14 genome groups from the full set of HBV genomes. In cases where the  
893 consensus length is less than the genome length, the given figure is the percentage of  
894 identical nucleotides (nts) in the matching region, not counting any alignment gaps or  
895 ambiguous consensus nts. For each sample, the genome group with the highest identity is  
896 highlighted in bold. **b**, Inter-consensus sequence identity. The Needleman-Wunsch algorithm  
897 was used to globally align all sample consensus sequences against one another. The table  
898 shows the nt identity percentage for each alignment. In cases where the consensus lengths  
899 were unequal, the given figure is the percentage of identical nts in the matching region, not  
900 counting any alignment gaps or ambiguous consensus nts.

901

#### 902 **Extended Data Figure 1 | Ancient DNA damage patterns**

903 The frequencies of the mismatches observed between the HBV reference sequences  
904 (Extended Data Table 1) and the reads are shown as a function of distance from the 5' end.  
905 C>T (5') and G>A (3') mutations are shown in red and blue, respectively. All other possible  
906 mismatches are reported in gray. Insertions are shown in purple, deletions in green, and  
907 clippings in orange. The count of reads matching HBV for each sample is shown in  
908 parentheses. **a**, Damage patterns for RISE563, DA222, DA119, RISE254, DA195, DA27,  
909 DA51, RISE386, RISE387, DA29, DA45, RISE154. **b**, Damage patterns for DA222 without  
910 (left) and with (right) USER treatment. **c**, Damage patterns with 10, 20, 50, 100, 200, 500  
911 and 1000 reads, where each opaque line corresponds to one replicate set of reads.

912

#### 913 **Extended Data Figure 2 | Hepadnavirus Maximum Likelihood tree**

914 Shows 26 sequences from the Orthohepadnavirus species (Dataset 1, see Methods)  
915 including the ancient HBV sequences. Ancient genotype A sequences are shown in red,  
916 ancient genotype B sequences in orange, ancient genotype D sequences in blue and novel  
917 genotype sequences in green. The tree was constructed in PhyML<sup>62</sup>, optimizing for topology,

918 branch lengths, and rates, with 100 bootstraps (see Methods). Internal nodes with <70%  
919 bootstrap support are shown as polytomies.

920

### 921 **Extended Data Figure 3 | Genotype A recombination break point evidence**

922 RDP4<sup>53</sup> was used to analyse the set of 12 ancient sequences plus a representative set of 15  
923 modern human and NHP sequences (see Methods). The seven recombination programs  
924 used by RDP4 suggested that all genotype A sequences are recombinants, with the  
925 genotype D sequence HBV-DA51 as the minor parent and an unknown major parent. The  
926 obvious interpretation is that recombination formed an ancestor of the oldest sequences,  
927 evidence of which is still present in the less ancient and the modern representatives. The  
928 panel shows the graphical evidence and predicted recombination break point distribution for  
929 the two oldest genotype A sequences, HBV-RISE386 and HBV-RISE387, according to three  
930 of the RDP4 methods (MaxChi, Bootscan, and RDP). In all sub-plots the predicted location  
931 of the break points is shown by a dashed vertical line and the surrounding gray area shows  
932 the 99% confidence interval for the break point. Sub-plots on the same row share their Y  
933 axis and those in the same column share their X axis. **a**, HBV-RISE386 analysed by  
934 MaxChi. **b**, HBV-RISE386 analysed by Bootscan. **c**, HBV-RISE386 analysed by RDP. **d**,  
935 HBV-RISE387 analysed by MaxChi. **e**, HBV-RISE387 analysed by Bootscan. **f**, HBV-  
936 RISE387 analysed by RDP.

937

### 938 **Extended Data Figure 4 | HBV Maximum likelihood tree**

939 The sequences from Dataset 2 (see Methods) and the ancient sequences were aligned in  
940 MAFFT<sup>61</sup>. The tree was constructed in PhyML<sup>62</sup>, optimizing for topology, branch lengths, and  
941 rates, with 100 bootstraps (see Methods). Internal nodes with <70% bootstrap support are  
942 shown as polytomies. Ancient genotype A sequences are shown in red, ancient genotype B  
943 sequences in orange, ancient genotype D sequences in blue and novel genotype sequences  
944 in green. Letters on internal branches indicate the genotype. Taxon names indicate:  
945 genotype / subgenotype, GenBank accession number, age, country abbreviation of

946 sequence origin, region of sequence origin, host species, and optional additional remarks.  
947 Note that the ML tree shows topological uncertainty (polytomies) in areas where the  
948 BEAST2<sup>26</sup> tree (Figure 2) is well resolved. This is the case for two reasons. Firstly, BEAST2  
949 always produces a fully-resolved binary topology without polytomies. Second, and more  
950 important, BEAST2 creates a time tree and uses tip dates to constrain the possible  
951 topologies under consideration. Thus BEAST2 can know that certain topologies are unlikely  
952 or impossible, whereas ML cannot and thus inherently has greater uncertainty regarding tree  
953 topology.

954

#### 955 **Extended Data Figure 5 | Root-to-tip regression and date randomisation tests**

956 **a**, Regression of root-to-tip distances and ages performed in Scipy<sup>65</sup>. 124 branch lengths  
957 were extracted using TempEst<sup>64</sup> from trees inferred using neighbour joining (NJ), ML, and  
958 Bayesian methods. Shaded areas show 95% confidence intervals. Slopes: 1.01E-05, 1.20E-  
959 05, 4.21E-06. Correlation coefficients: 0.45 (R<sup>2</sup>=0.2), 0.36 (R<sup>2</sup>=0.13), 0.51 (R<sup>2</sup>=0.26) for  
960 ML, Bayesian, and NJ trees, respectively. **b**, Date randomisation tests under the strict clock  
961 model. The median and 95% HPD interval for the substitution rates are given. The rate for  
962 the correctly dated tree is shown in red. Dates were randomised within all sequences, within  
963 the ancient sequences only, and within each genotype. We performed three replicates of  
964 each. None of the 95% HPD intervals for the randomised runs overlap with the 95% HPD  
965 intervals for the correctly dated runs, suggesting the presence of a temporal signal in the  
966 data.

967