# A 20+ Ma old enamel proteome from Canada's High Arctic reveals diversification of Rhinocerotidae in the middle Eocene-Oligocene

## Authors

Ryan S. Paterson[1*], Meaghan Mackie[1,2], Alessio Capobianco[3,4], Nicola S. Heckeberg[3,4], Danielle Fraser[5,6,7,8], Fazeelah Munir[9], Ioannis Patramanis[1], Jazmín Ramos-Madrigal[1], Shanlin Liu[10], Abigail D. Ramsøe[1], Marc R. Dickinson[9], Chloë Baldreki[9], Marisa Gilbert[5], Raffaele Sardella[11], Luca Bellucci[12], Gabriele Scorrano[1,13], Fernando Racimo[1], Eske Willerslev[1,14,15], Kirsty E.H. Penkman[9], Jesper V. Olsen[2], Ross D.E. MacPhee[16], Natalia Rybczynski[5*], Sebastian Höhna[3,4], Enrico Cappellini[1*]

* Corresponding authors: Ryan Sinclair Paterson: ryan.paterson@sund.ku.dk; Natalia Rybczynski: nrybczynski@nature.ca; Enrico Cappellini: ecappellini@sund.ku.dk

## Affiliations

[1]Globe Institute, University of Copenhagen, 5-7 Øster Voldgade, 1350, Copenhagen, Denmark, [2]Novo Nordisk Foundation Center for Protein Research, University of Copenhagen, Blegdamsvej 3b, 2200, Copenhagen, Denmark, [3]GeoBio-Center LMU, Ludwig-Maximilians-Universität München, Richard-Wagner-Str. 10, Munich 80333, Germany, [4]Department of Earth and Environmental Sciences, Palaeontology & Geobiology, Ludwig-Maximilians-Universität München, Richard-Wagner-Str. 10, Munich 80333, Germany, [5]Palaeobiology, Canadian Museum of Nature, PO Box 3443 Stn "D", Ottawa ON K1P 6P4, Canada, [6]Department of Earth Sciences, Carleton University, 1125 Colonel By Drive, Ottawa, Ontario K1S 5B6, Canada, [7]Department of Biology, Carleton University, 1125 Colonel By Drive, Ottawa, Ontario K1S 5B6, Canada, [8]Department of Paleobiology, Smithsonian Institution, National Museum of Natural History, 10th and Constitution NW, Washington, DC 20013-

23  7012, [9]Department of Chemistry, University of York, York, YO10 5DD, UK, [10]Department of

24  Entomology, College of Plant Protection, China Agricultural University, Beijing, China, [11]Sapienza

25  University of Rome, Department of Earth Sciences (PaleoFactory lab.), Piazzale Aldo Moro 5, 00185,

26  Rome, Italy, [12]Museo di Geologia e Paleontologia, Sistema Museale di Ateneo, Università di Firenze,

27  via Giorgio La Pira 4, 50121, Firenze, Italy, [13]Center for Molecular Anthropology for the Study of

28  Ancient DNA, Department of Biology, University of Rome "Tor Vegata", via della Ricerca

29  Scientifica n. 1, 173, Rome, Italy, [14]Department of Zoology, University of Cambridge, 17 Mill Lane,

30  CB2 1RX, Cambridge, UK, [15]Marum Center for Marine Environmental Sciences and Faculty of

31  Geosciences, University of Bremen, Leobener Str. 8, 28359, Bremen, Germany, [16]Department of

32  Mammalogy, American Museum of Natural History, 200 Central Park West, 10024, New York City,

33  United States.

34

35  **In the past decade, ancient protein sequences have emerged as a valuable source of data**

36  **for deep-time phylogenetic inference. Still, the recovery of protein sequences providing novel**

37  **phylogenetic insights does not exceed 3.7 Ma (Pliocene). Here, we push this boundary back to**

38  **21-24 Ma (early Miocene), by retrieving enamel protein sequences of an early-diverging**

39  **rhinocerotid (*Epiaceratherium* sp. - CMNF-59632) from the Canadian High Arctic. We recover**

40  **partial sequences of seven enamel proteins (AHSG, ALB, AMBN, AMELX, AMTN, ENAM,**

41  **MMP20) and over 1000 peptide-spectrum matches, spanning over at least 251 amino acids.**

42  **Authentic endogeneity of these sequences is supported by indicators of protein damage,**

43  **including several spontaneous and irreversible post-translational modifications accumulated**

44  **during prolonged diagenesis and reaching near-complete occupancy at many sites. Bayesian tip-**

45  **dating, across 15 extant and extinct perissodactyl taxa, places the divergence time of CMNF-**

46  **59632 in the middle Eocene-Oligocene, and identifies a later divergence time for**

47  **Elasmotheriinae in the Oligocene. The finding weakens alternative models suggesting a deep**

48  **basal split between Elasmotheriinae and Rhinocerotinae. This divergence time of CMNF-59632**

49  **coincides with a phase of high diversification of rhinocerotids, and supports a Eurasian origin of**

50    **this clade in the late Eocene or Oligocene. The findings are consistent with previous hypotheses**

51    **on the origin of the enigmatic fauna of the Haughton crater, which, in spite of their considerable**

52    **degree of endemism, also display similarity to distant Eurasian faunas. Our findings**

53    **demonstrate the potential of palaeoproteomics in obtaining phylogenetic information from a**

54    **specimen that is ten times older than any sample from which endogenous DNA has been**

55    **obtained.**

56    Phylogenetic placement of deep-time (>1 Ma) fossils has typically relied on morphological

57    observations, as the recovery of sufficiently extensive genetic evidence has not been thought to be

58    possible beyond the Pleistocene[1]. While ancient DNA (aDNA) sequences are a valuable source of

59    data for inferring phylogenies and population dynamics in the late Pleistocene[2–5], the oldest authentic

60    aDNA from macrofossils has been extracted from Arctic-situated specimens dated to no more than 1.2

61    Ma[6]. In contrast, palaeoproteomic data have been recovered from late Miocene, Pliocene and early

62    Pleistocene fossils, even in localities that are warm, humid, and/or at low latitudes[7,8]. Although protein

63    sequences from the early Pleistocene have been successfully used to infer the phylogenetic placement

64    of various fossil mammals[9–11], the precise limit of proteomic survival has not yet been systematically

65    characterised[12,13]. Currently, the oldest confirmed palaeoproteomic data successfully used to infer sub-

66    ordinal taxonomic relationships derive from bone collagen of camelids from the 3.7 Ma Fyles Leaf

67    Bed site of Canada's High Arctic[14,15]. This illustrates how our understanding of evolutionary

68    relationships is currently limited by the preservation of biomolecules from extinct species.

69    Rhinocerotidae is a family that includes only five extant species, but a wide diversity of fossil

70    members[16,17]. It remains debated as to where and when the radiation of this group occurred[18]. For most

71    of the last two decades, the group was defined by a deep 'basal split' between two clades -

72    Rhinocerotinae and Elasmotheriinae - prior to bursts of rhinocerotid diversification in the late

73    Eocene[19–23]. This paradigm contrasts earlier hypotheses of a close relationship between two extinct

74    rhinocerotids that survived into the late Pleistocene – the Siberian unicorn (Elasmotheriinae,

75    *Elasmotherium sibiricum*) and the woolly rhino (Rhinocerotinae, *Coelodonta antiquitatis*)[24]. Recently,

76    the sequenced genomes of *Coelodonta* and *Elasmotherium*[25] were used to confirm hypotheses based

77      on morphological data that suggest they have distinct phylogenetic affinities[19], but also allowed for

78      the recognition of a relatively-young split between these two groups during the late Eocene (36 Ma).

79      This suggests that the deep-divergence hypothesis based on the morphological analysis of fossils is

80      not supported by molecular evidence. However, the lack of available genetic sequence data from other

81      early-diverging rhinocerotid lineages makes it difficult to assess the timing of the Rhinocerotinae-

82      Elasmotheriinae split in relation to other radiations that occurred within the group. For these reasons,

83      the ancient radiations of the group still remain obscured.

84      To investigate the timing of Rhinocerotidae divergence and the potential for evolutionarily

85      informative protein sequences to persist in deep-time, we targeted vertebrate dental enamel deriving

86      from the Haughton crater (75°N, Nunavut) in Canada's High Arctic (Figure 1). The Haughton crater

87      is an impact structure with its stratigraphy including post-impact fossiliferous lacustrine sediments

88      dated to 21-24 Ma[26]. While their geological age is advanced, fossils from these sediments are found in

89      a polar landscape, currently characterised by permafrost. This creates a temperature regime favourable

90      for biomolecular preservation, sparing these fossils from the harshest effects of diagenesis.

91      Accordingly, specimens from this site serve as promising candidates for biomolecular preservation in

92      deep-time.

93      The digestion-free palaeoproteomic workflow[9,10] applied to an early Miocene rhinocerotid

94      (*Epiaceratherium* sp.) specimen of dental enamel[27] from the Haughton Formation (21.8 Ma) allowed

95      for the recovery of an enamel proteome covering 1163 confident peptide-spectrum matches (PSMs),

96      at least seven proteins (AHSG, ALB, AMBN, AMELX, AMTN, ENAM, MMP20), and spanning at

97      least 251 amino acids (Figure 2a) (Ext. Figure 1a). The enamel proteome of CMNF-59632 currently

98      represents both the oldest mammalian skeletal proteome currently reported, confirming the predicted

99      deep-time persistence of ancient mammalian proteins from high latitudes[8,10], and the first

100      biomolecular characterisation of the extinct genus *Epiaceratherium*. While the survival of a relatively

101      rich enamel proteome from such ancient deposits is surprising, the age of the specimen belies its

102      excellent state of preservation.

103    To better appreciate the preservation state of the Haughton Crater enamel proteome, we

104    compared it to those of two other rhinocerotids, namely the early Pleistocene *Stephanorhinus* from the

105    site of Dmanisi (Georgia), dated at 1.77 Ma[9], and a middle Pleistocene *Stephanorhinus* (~0.4 Ma)

106    from the site of Fontana Ranuccio (Italy). While the set of proteins retrieved from the CMNF-59632

107    enamel specimen is similar to that of the other two Pleistocene rhinocerotids used for comparison

108    (Figure 2a), fewer peptides and a shorter reconstructed amino acid sequence were recovered from the

109    Arctic specimen. In addition to proteins previously found in deep-time enamel samples (AHSG,

110    AMBN, AMELX, AMTN, ENAM, MMP20), serum albumin was also found in each rhinocerotid

111    sample. Previously recovered from other fossil enamel specimens[11], serum albumin is

112    phylogenetically more informative than enamel-specific proteins due to its higher amino acid

113    sequence variability[28]. In the present analysis, while some albumin peptides are removed during initial

114    data filtering (as they are identical to peptide sequences from bovine serum albumin, a potential

115    laboratory contaminant), several others do not match contaminants. The vast majority of the spectra

116    that confidently support their identification show post-translational modifications (PTMs) that derive

117    from prolonged diagenesis, which supports their authenticity.

118    As expected, diagenetic modifications and PTMs are extensive in the enamel proteome of

119    CMNF-59632 (Figure 2b). Average peptide lengths are similar, though slightly shorter than those of

120    the Dmanisi early Pleistocene specimen, and further reduced in comparison to the Fontana Ranuccio

121    middle Pleistocene *Stephanorhinus*, indicating a greater degree of peptide bond hydrolysis (Figure

122    2b). We also observe high deamidation rates in CMNF-59632, though no more so than in the

123    Pleistocene rhinocerotids (Ext. Figure 2). While high deamidation rates can be useful for confirming

124    proteome authenticity, they can be highly variable within samples[29–31], and can plateau relatively-

125    quickly in fossil proteomes, reducing their utility in characterising degradation patterns in deep-time

126    (Figure 2c). Instead, we identify a suite of informative spontaneous PTMs indicative of advanced

127    diagenesis that are observed at a higher rate in the Arctic Miocene rhinocerotid, providing support for

128    their utility as markers of advanced diagenesis and authenticity in deep-time (Figure 2c)[9]. These

129    include arginine to ornithine conversion (Figure 2c) and advanced forms of tryptophan (Ext. Figure

130    2b) and histidine oxidation (Ext. Figure 1c). Intra-crystalline protein decomposition analysis further

131    confirms the advanced degradation state of CMNF-59632. The concentration of free and total

132    hydrolysable amino acids (FAA and THAA, respectively) is around half of those in the early

133    Pleistocene *Stephanorhinus* sample from Dmanisi (Extended Figure 2a), and the percentage of FAA

134    in CMNF-59632 (~75%) is higher than in the Pleistocene rhino from Dmanisi (~50%) (Ext. Figure

135    3b), supporting increased peptide bond hydrolysis. Furthermore, the racemisation values for CMNF-

136    59632 fall along the expected FAA vs THAA trends for both fossil enamel, and the experimentally-

137    heated samples (Ext. Figure 4), confirming the closed system behaviour of CMNF-59632 enamel

138    amino acids, and supporting the endogeneity of the peptides retrieved.

139          At least 10 single amino acid polymorphisms (SAPs) support the placement of CMNF-59632

140    within Rhinocerotidae. A smaller number (2+) of SAPs are shared between CMNF-59632 and other

141    perissodactyls, to the exclusion of later-diverging rhinocerotids. No novel variants are uncovered in

142    CMNF-59632, as the aforementioned SAPs represent character states retained from ancestors within

143    Perissodactyla and Mammalia more broadly. The identification of these SAPs is supported by several

144    unique PSMs displaying almost complete ion series (e.g., Ext. Figure 5).

145          Peptide sequences recovered from CMNF-59632 derive from similar sequence regions to

146    those previously identified in the Dmanisi Pleistocene *Stephanorhinus* proteome (Figure 2d),

147    particularly for the three most abundant enamel matrix proteins (EMPs). ENAM and AMBN present

148    broadly similar sequence coverage patterns in both specimens, though with fewer PSMs covering

149    most positions in the Miocene sample. AMELX, the most abundant EMP, is instead covered by a

150    similar number of PSMs in both the Miocene and Pleistocene samples. The depth of coverage is also

151    similar for the most abundantly-covered AMELX sequences, including those spanning the deletion

152    observed in the Leucine-Rich Amelogenin Peptide (LRAP)[9,32].

153          Regardless of the mechanisms behind preferential mass spectrometric and data analysis

154    identification of specific sequence regions, biases favouring the recovery[7] and identification[33] of

155    conserved peptide sequences can ultimately lead to underestimates of divergence times in taxa

156    represented by empirically-derived protein sequences. To accurately estimate the phylogenetic

157    position of CMNF-59632 and estimate divergence times within the group, we completed a

158    phylogenetic analysis of a suite of extinct and extant perissodactyls. In addition to the perissodactyl

159    taxa used in Cappellini et al. (2019)[9], we incorporated whole-genome sequence data to predict enamel

160    protein sequences from the Siberian unicorn (*Elasmotherium sibiricum*) and a pair of extant tapirs

161    (*Tapirus terrestris* and *Tapirus indicus*).

162        The time-calibrated phylogenetic analysis of enamel protein sequences under a Fossilised Birth

163    Death (FBD) model infers CMNF-59632 as the earliest diverging rhinocerotid in the analysis, with

164    *Elasmotherium sibiricum* being more closely related to Rhinocerotina (crown rhinoceroses) than to

165    CMNF-59632 (Figure 4). This phylogenetic hypothesis is also supported by Fraser et al. (2024) in a

166    total-evidence analysis. Additionally, our FBD analysis resolves the early Pleistocene *Stephanorhinus*

167    from Dmanisi as a sampled ancestor of the middle Pleistocene *Stephanorhinus* from Fontana Ranuccio.

168    Divergence time estimates place the split between CMNF-59632 and all other rhinocerotids during the

169    middle Eocene–Oligocene (around 41–25 Ma). The divergence between *Elasmotherium sibiricum* and

170    Rhinocerotina is reconstructed to have likely occurred in the Oligocene (around 34–22 Ma), which is

171    younger than previous molecular clock estimates[25].

172        The late Eocene and the early Oligocene represent dynamic periods in the evolution of

173    rhinocerotids, particularly in North America. After appearing in the middle Eocene (37-34 Ma)[34], North

174    American rhinocerotids diversify during the late Eocene, evolving a variety of body sizes and ecologies

175    as several new clades arise, before rhinocerotid diversity experiences a significant drop in the early

176    Oligocene (34-32 Ma)[35]. During this timeframe, other early-diverging lineages are also appearing in

177    Asia[18,36,37], and eventually spreading as far as Western Europe[18]. Morphologically, the Haughton crater

178    rhinocerotid shares closer affinities with these early-diverging lineages from Eurasia[38], particularly

179    those within the genus *Epiaceratherium*[27]. Similarly, some other vertebrates within the highly-endemic

180    fauna of the Haughton Formation have their closest relatives in Eurasia. These include the transitional

181    pinniped *Puijila darwini*, sister to the Oligocene *Potamotherium* of Europe[39], and the swan-like anatid,

182    a group which is otherwise restricted to the Oligocene and Miocene of Europe[38]. Overall, these patterns,

183    in conjunction with the recovered divergence times, suggest the Haughton crater rhinocerotid represents

184    a migrant from eastern Asia or western Europe, derived from one of the early-diverging lineages that

185    arose in the late Eocene or early Oligocene of East Asia.

186        We provide molecular evidence that this lineage falls outside of Rhinocerotinae, as it diverges

187    before the Rhinocerotinae-Elasmotheriinae split. We also reject a deep-divergence (basal split) between

188    Elasmotheriinae and Rhinocerotinae[19–23] and find moderate support for their branching event after the

189    divergence of *Epiaceratherium*. Our analysis disagrees with that of Kosintsev et al. (2019)[23], who find

190    a deep divergence for Elasmotheriinae (47.3 Ma), and an early divergence for Rhinocerotinae (almost

191    30.8 Ma). The later divergence times for these nodes in our analysis are in spite of equivalently old ages

192    for crown Ceratomorpha (earliest Eocene). Among other timetrees, our dates are generally most

193    consistent with those of Liu et al. (2021)[25]. Our recovered topologies are also broadly similar to trees

194    derived from the morphology-based phylogenetic analyses of Tissier et al. (2020)[18] and Lu et al.

195    (2023)[40], identifying Elasmotheriinae and Rhinocerotinae as deeply-nested within Rhinocerotidae.

196    Discrepancies between the genomic [25] and proteomic trees arise likely due to different calibration

197    points. The more ancient age of Elasmotheriinae in the analysis of Liu et al. (2021)[25] is constrained by

198    a high minimum bound for the Elasmotheriinae-Rhinocerotinae split (35 Ma). However, this date is

199    based on the earliest age of *Epiaceratherium naduongense* and its allocation to Rhinocerotinae.

200    Assuming monophyly of *Epiaceratherium*, the present proteomic evidence refutes the assignment of

201    this genus to Rhinocerotinae, as it falls as earlier-diverging than Elasmotheriinae without such

202    topological constraints in our phylogenetic analysis.

203        In sum, these findings highlight the importance of integrating palaeoproteomic sequence data

204    into phylogenetic analyses to infer topologies and estimate divergence times. Ancient proteomic

205    sequence data allows for robustly-supported timetrees, and can serve to develop phylogenetic

206    frameworks in deep-time, particularly from specimens too old to preserve ancient DNA. For example,

207    the present data allows for firm placement of the Haughton Crater rhinocerotid outside of Rhinocerotina,

208    and likely outside the Elasmotheriinae-Rhinocerotinae clade, a fact which has significant implications

209    for both morphological and molecular studies integrating fossil calibration times from the fossil record.

210 Furthermore, we demonstrate the deep-time survival of a rich set of peptides derived from proteins

211 present in mammalian enamel, well beyond the previously known limits of survival. This work

212 illustrates the power of palaeoproteomics in elucidating phylogeny and taxonomy of extinct vertebrates

213 in deep-time. These findings should encourage further vertebrate palaeontological fieldwork in the High

214 Arctic, and other cold-temperature sites with taphonomic conditions favourable to biomolecular

215 preservation.

## References

217 1. Kjær, K. H. *et al.* A 2-million-year-old ecosystem in Greenland uncovered by environmental

218 DNA. *Nature* **612**, 283–291 (2022).

219 2. Orlando, L. *et al.* Recalibrating Equus evolution using the genome sequence of an early Middle

220 Pleistocene horse. *Nature* **499**, 74–78 (2013).

221 3. Dabney, J. *et al.* Complete mitochondrial genome sequence of a Middle Pleistocene cave bear

222 reconstructed from ultrashort DNA fragments. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 15758–15763

223 (2013).

224 4. Meyer, M. *et al.* A mitochondrial genome sequence of a hominin from Sima de los Huesos.

225 *Nature* **505**, 403–406 (2014).

226 5. Barlow, A. *et al.* Middle Pleistocene genome calibrates a revised evolutionary history of extinct

227 cave bears. *Curr. Biol.* **31**, 1771–1779.e7 (2021).

228 6. van der Valk, T. *et al.* Million-year-old DNA sheds light on the genomic history of mammoths.

229 *Nature* **591**, 265–269 (2021).

230 7. Demarchi, B. *et al.* Protein sequences bound to mineral surfaces persist into deep time. *Elife* **5**,

231 (2016).

232 8. Demarchi, B. *et al.* Survival of mineral-bound peptides into the Miocene. *Elife* **11**, (2022).

233 9. Cappellini, E. *et al.* Early Pleistocene enamel proteome from Dmanisi resolves Stephanorhinus

234 phylogeny. *Nature* **574**, 103–107 (2019).

235 10. Welker, F. *et al.* Enamel proteome shows that Gigantopithecus was an early diverging pongine.

236     *Nature* **576**, 262–265 (2019).

237     11.  Welker, F. *et al.* The dental proteome of Homo antecessor. *Nature* **580**, 235–238 (2020).

238     12.  Buckley, M., Warwood, S., van Dongen, B., Kitchener, A. C. & Manning, P. L. A fossil protein

239          chimera; difficulties in discriminating dinosaur peptide sequences from modern cross-

240          contamination. *Proc. Biol. Sci.* **284**, (2017).

241     13.  Umamaheswaran, R. & Dutta, S. Preservation of proteins in the geosphere. *Nat Ecol Evol* (2024)

242          doi:10.1038/s41559-024-02366-z.

243     14.  Rybczynski, N. *et al.* Mid-Pliocene warm-period deposits in the High Arctic yield insight into

244          camel evolution. *Nat. Commun.* **4**, 1550 (2013).

245     15.  Buckley, M., Lawless, C. & Rybczynski, N. Collagen sequence analysis of fossil camels,

246          Camelops and c.f. Paracamelus, from the Arctic and sub-Arctic of Plio-Pleistocene North

247          America. *J. Proteomics* **194**, 218–225 (2019).

248     16.  Prothero, D. R. Rhinocerotidae. in *Evolution of tertiary mammals of North America* 595 (1998).

249     17.  Pandolfi, L. Evolutionary history of Rhinocerotina (Mammalia, Perissodactyla). *Fossilia* **2018**,

250          27–32 (2018).

251     18.  Tissier, J., Antoine, P.-O. & Becker, D. New material of Epiaceratherium and a new species of

252          Mesaceratherium clear up the phylogeny of early Rhinocerotidae (Perissodactyla). *R Soc Open*

253          *Sci* **7**, 200633 (2020).

254     19.  Antoine, P. O. & national d'histoire naturelle, M. Phylogénie et évolution des Elasmotheriina

255          (Mammalia, Rhinocerotidae). (2002).

256     20.  Antoine, P.-O. Middle Miocene elasmotheriine Rhinocerotidae from China and Mongolia:

257          taxonomic revision and phylogenetic relationships. *Zool. Scr.* **32**, 95–118 (2003).

258     21.  Antoine, P.-O. *et al.* A revision of Aceratherium blanfordi Lydekker, 1884 (Mammalia:

259          Rhinocerotidae) from the Early Miocene of Pakistan: postcranials as a key. *Zool. J. Linn. Soc.*

260          **160**, 139–194 (2010).

261     22.  Becker, D., Antoine, P.-O. & Maridet, O. A new genus of Rhinocerotidae (Mammalia,

262          Perissodactyla) from the Oligocene of Europe. *J. Syst. Palaeontol.* **11**, 947–972 (2013).

263     23.  Kosintsev, P. *et al.* Evolution and extinction of the giant rhinoceros Elasmotherium sibiricum

264   sheds light on late Quaternary megafaunal extinctions. *Nat Ecol Evol* **3**, 31–38 (2019).

265 24. Cerdeño, E. *Cladistic Analysis of the Family Rhinocerotidae (Perissodactyla)*. (American

266   Museum of Natural History, 1995).

267 25. Liu, S. *et al.* Ancient and modern genomes unravel the evolutionary history of the rhinoceros

268   family. *Cell* **184**, 4874–4885.e16 (2021).

269 26. Jessberger, E. K. 40Ar-39Ar Dating of the Haughton Impact Structure. *Meteoritics* (1988).

270 27. Fraser D, Rybczynski N, Gilbert M, Dawson MR. Post-Eocene Rhinocerotid Dispersal via the

271   North Atlantic. *BioRxiv* (2024) doi:10.1101/2024.06.04.597351.

272 28. Nei, M. *Molecular Evolutionary Genetics*. (Columbia University Press, 1987).

273 29. Simpson, J. P. *et al.* The effects of demineralisation and sampling point variability on the

274   measurement of glutamine deamidation in type I collagen extracted from bone. *J. Archaeol. Sci.*

275   **69**, 29–38 (2016).

276 30. Brown, S. *et al.* Examining collagen preservation through glutamine deamidation at Denisova

277   Cave. *J. Archaeol. Sci.* **133**, 105454 (2021).

278 31. Pal Chowdhury, M. & Buckley, M. Trends in deamidation across archaeological bones, ceramics

279   and dental calculus. *Methods* **200**, 67–79 (2022).

280 32. Gibson, C. W. *et al.* Identification of the leucine-rich amelogenin peptide (LRAP) as the

281   translation product of an alternatively spliced transcript. *Biochem. Biophys. Res. Commun.* **174**,

282   1306–1312 (1991).

283 33. Welker, F. Elucidation of cross-species proteomic effects in human and hominin bone proteome

284   identification through a bioinformatics experiment. *BMC Evol. Biol.* **18**, 23 (2018).

285 34. Hanson, C. B. Teletaceras Radinskyi, a new primitive Rhinocerotid from the Late Eocene Clarno

286   Formation of Oregon. in *The evolution of perissodactyls* 379–398 (1989).

287 35. Prothero, D. R. *The Evolution of North American Rhinoceroses*. (Cambridge University Press,

288   2005).

289 36. Antoine, P.-O. *et al.* Early rhinocerotids (Mammalia: Perissodactyla) from South Asia and a

290   review of the Holarctic Paleogene rhinocerotid record. *Can. J. Earth Sci.* **40**, 365–374 (2003).

291 37. Böhme, M. *et al.* Na Duong (northern Vietnam)-an exceptional window into Eocene ecosystems
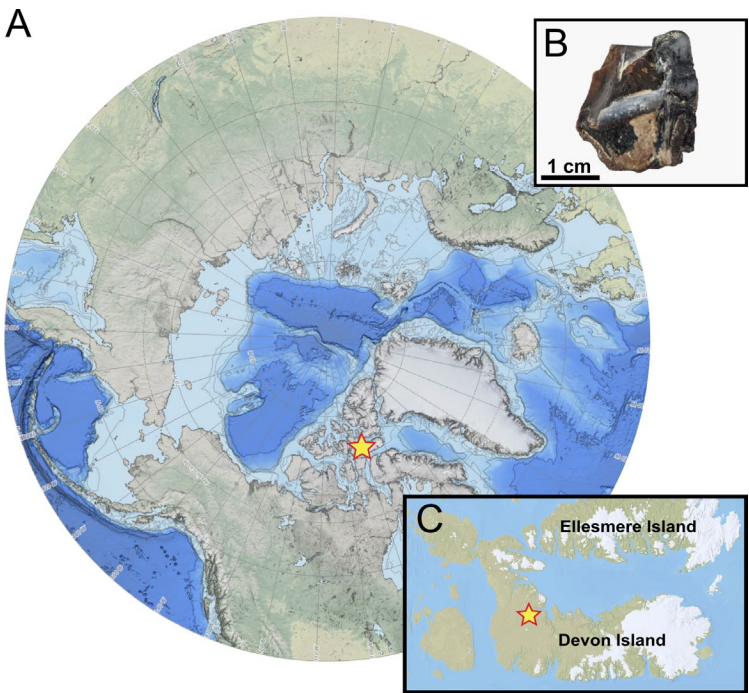
292  from Southeast Asia. *Zitteliana. Reihe A, Mitteilungen der Bayerischen Staatssammlung für*

293  *Paläontologie und Geologie* **53**, 121–167 (2013).

294 38. Whitlock, C. & Dawson, M. R. Pollen and Vertebrates of the Early Neogene Haughton

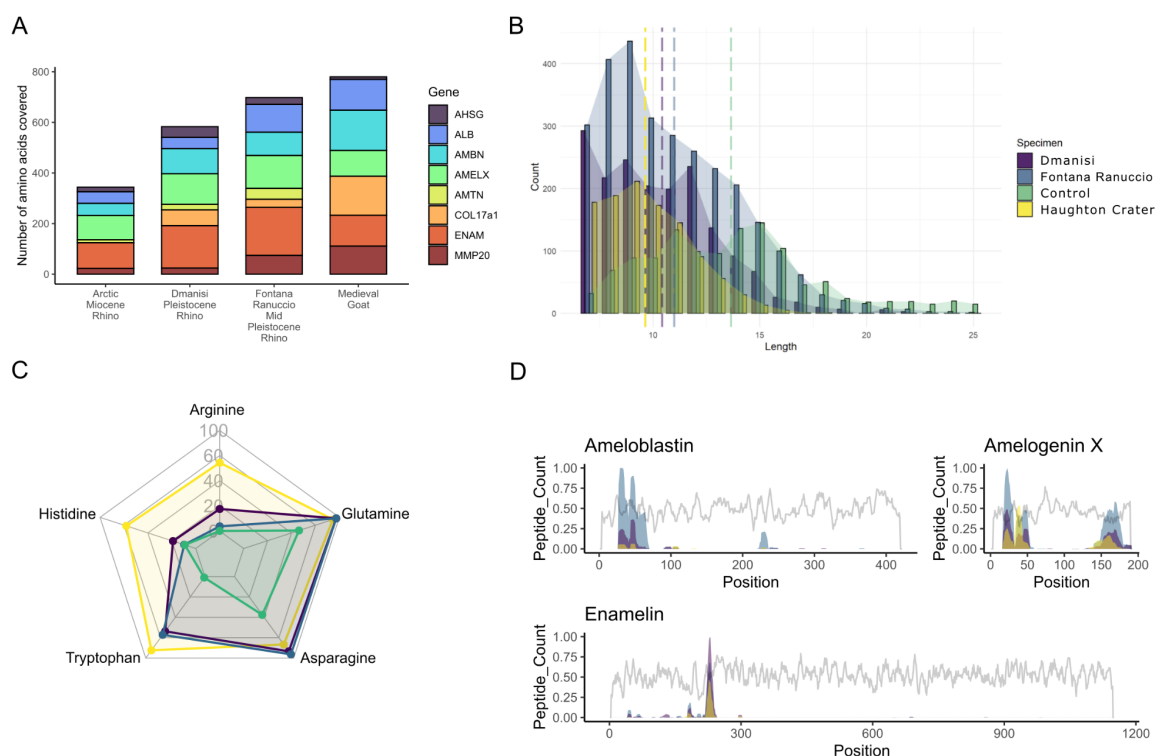295  Formation, Devon Island, Arctic Canada. *Arctic* **43**, 324–330 (1990).

296 39. Rybczynski, N., Dawson, M. R. & Tedford, R. H. A semi-aquatic Arctic mammalian carnivore

297  from the Miocene epoch and origin of Pinnipedia. *Nature* **458**, 1021–1024 (2009).

298 40. Lu, X.-K., Deng, T. & Pandolfi, L. Reconstructing the phylogeny of the hornless rhinoceros

299  Aceratheriinae. *Frontiers in Ecology and Evolution* **11**, (2023).

300

301

302

303

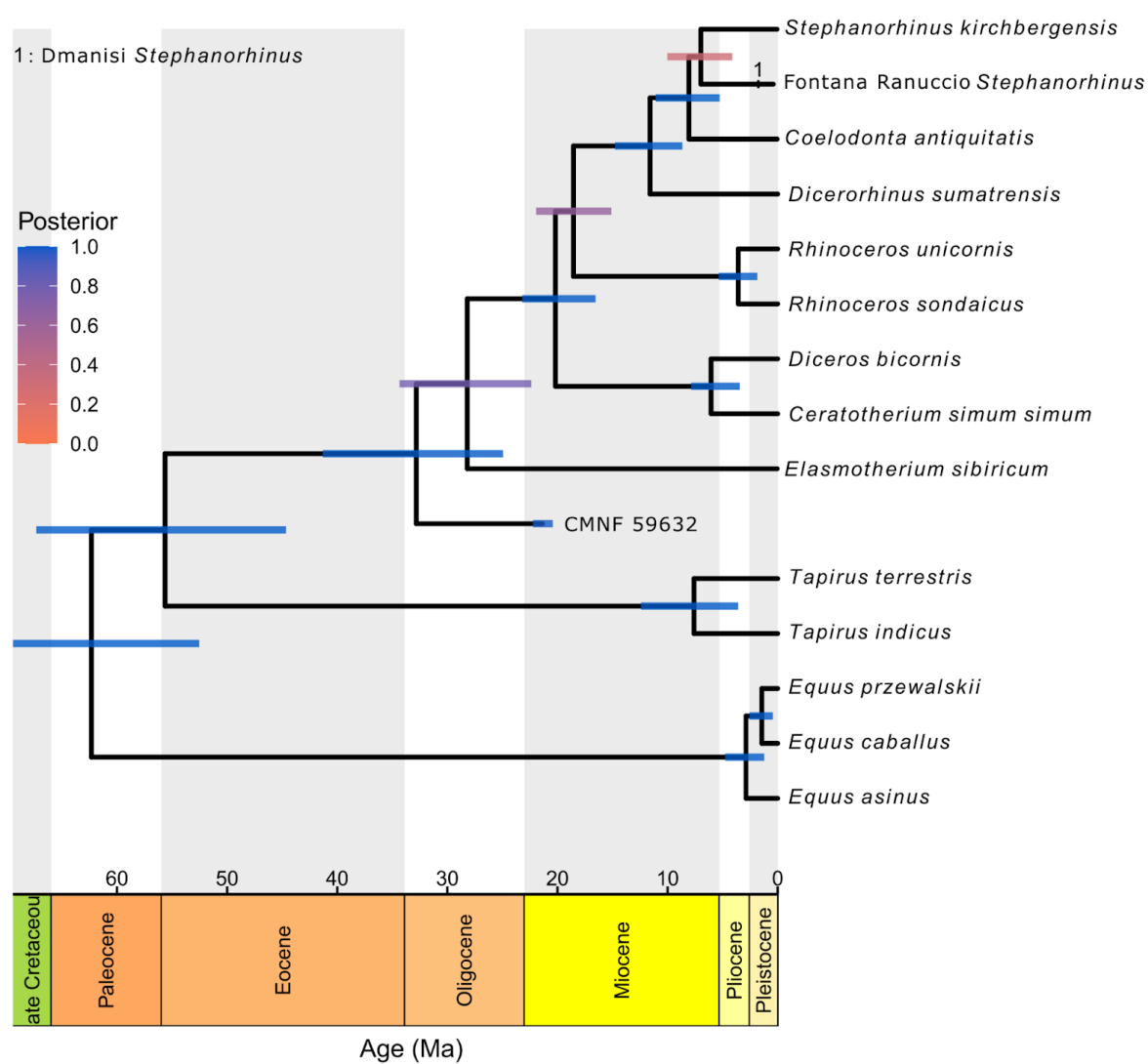304

305

306

307

308

309

310

Figures



**Figure 1) The high latitude Haughton crater on Devon Island has produced a highly endemic vertebrate fauna.** A) Location of Devon Island within the circumpolar North. B) Anterolingual view of specimen CMNF-59632 after destructive palaeoproteomic analysis. C) Location of Haughton crater on Devon Island.

**Figure 2) Proteome preservation in the enamel specimen of the early Miocene rhinocerotid (CMNF-59632).** Preservation is compared to enamel proteomes from an early Pleistocene (1.77 Ma) *Stephanorhinus* (DM.5/157), a middle Pleistocene (0.4 Ma) *Stephanorhinus* (CGG 1_023342) and a medieval ovicaprine (Control)[9]. All plots exclude contaminants and reverse hits. A) Amino acid sequence coverage for each identified protein; B) peptide length distributions, dashed bars represent average peptide length for each specimen; C) proportion of a selected sample of amino acids that are often modified in ancient enamel proteomes. Results derive from PTM-specific searches described in methods. 'Arginine' includes arginine-to-ornithine conversion, 'Glutamine' includes glutamine deamidation, 'Asparagine' includes asparagine deamidation, 'Tryptophan' includes advanced tryptophan oxidation to kynurenine, oxylactone, and tryptophandione, 'Histidine' includes oxidation and dioxidation of histidine, histidine conversion to hydroxyglutamate. D) Sequence coverage plots for the three most abundant enamel matrix proteins (AMBN, AMELX, ENAM), recording number of PSMs (coloured areas) and mutability (grey line).

**Figure 3) Abridged alignment and mirror plots of a phylogenetically-informative single amino acid polymorphism (SAP) at AMELX-39.** The top spectrum is experimentally-derived, while the bottom one is predicted using the 'Original mode' with the Prosit tool, available online via the Universal Spectrum Explorer[83]. This spectrum is the highest scoring peptide-spectrum match (with Andromeda) for AMELX sequence positions spanning the most abundantly-covered SAP differentiating between CMNF-59632, and all other rhinocerotids for which sequences are available.

**Figure 4) Time-calibrated phylogeny of Rhinocerotidae enamel proteomes.** The maximum *a posteriori* (MAP) tree was produced using RevBayes v.1.2.1 [70]; https://revbayes.github.io, with a Fossilized Birth Death (FBD) model. Coloured bars at nodes represent 95% height posterior density (HPD) age interval estimates. Specimen CMNF-59632 represents the early Miocene rhinocerotid from the Haughton Crater.

## Materials and Methods

### Site and Specimen

Located within the Haughton impact crater (75°N, Nunavut, Canada), the Haughton Formation comprises the remnants of a large, post-impact lacustrine deposit, dated to the early Miocene. Previous dating estimates, using fission-track and $^{40}$Ar–$^{39}$Ar furnace step-heating dating, identified an age of 24-21 Ma[26,41]. An early Miocene age has also been corroborated by (U-Th)/He thermochronology[42]. While older age estimates between 30-40 Ma have also been suggested[43–45], there have been no age estimates younger than the early Miocene. Therefore, we conservatively use the younger early Miocene age estimates in our analysis and interpretation.

The highly-endemic fauna of the Haughton formation consists of several vertebrate taxa, including a transitional pinniped[39], a pair of salmoniform fishes, a swan-like anatid, a small artiodactyl, a leporid rabbit, a heterosocid shrew, and a well-preserved rhinocerotid[27,38]. While the megafloral assemblage is not particularly rich, the palynofloral assemblage is well-characterised, allowing for reconstruction of local climatic conditions. In the early Miocene, the Haughton crater lake and its surrounding environs experienced a significantly warmer annual temperature (8-12°C) than the present day[38,46].

The specimen CMNF-59632 is a nearly complete rhinocerotid skeleton, including skull and dentition,  uncovered 10.8 m above the base of the formation[27]. The present analysis focuses on a single tooth fragment from a lower left m1 (Figure 1b) that was already separated from the rest of its tooth row due the fragmenting effectings of cryoturbation. The dental specimen's rhinocerotid affinities are further supported by its size and morphology, most notably the presence of vertical Hunter-Schreger bands on its enamel, a defining feature of rhinocerotids and found in few other mammals[47]. A single tusk fragment (left i2) derived from CMNF-59632 was also selected for proteomic extraction. Due to its thin enamel, only limited peptides were recovered from this tusk fragment, and the sample is thus excluded from further analysis and discussion.

376    Proteomic extraction and LC-MS/MS

377        The laboratory workflow for the CMNF-59632 teeth and the Fontana Ranuccio

378    *Stephanorhinus* tooth (for comparison) generally follows that of Cappellini et al. (2019) and Taurozzi

379    et al. (2024)[9,48]. Using a sterilised drill, flakes of enamel were removed from the fragmentary teeth,

380    with care taken to avoid sampling the dentine. The CMNF-59632 tooth enamel sample, 154 mg, was

381    then ground to a fine powder, and demineralized overnight using 10% HPLC-grade trifluoroacetic

382    acid (TFA) (Merck, Sigma-Aldrich). The CMNF-59632 tusk enamel sample, weighing 90 mg, was

383    processed in the same way. The Fontana Ranuccio enamel sample was divided into three subsamples -

384    FR2, FR3, and FR4 - weighing 202, 243, and 205 mg, respectively. They were similarly ground to a

385    fine powder, and demineralized using 10% TFA (FR3, FR4) or 10% HCl (hydrochloric acid, FR2).

386    For each sample, the demineralization step was repeated a second time to ensure complete

387    demineralization. As enamel peptides are already hydrolysed *in vivo*, no enzymatic digestion was

388    performed. Subsequently, peptides were collected and desalted on C-18 StageTips[49] produced in-

389    house. An extraction blank for each sample set was processed alongside the samples for every step, to

390    control for contamination.

391    Mass spectrometry

392        Stagetips were eluted with 30 μL of 40% acetonitrile (ACN) and 0.1% formic acid (FA) into a

393    96-well plate. To remove ACN and concentrate the samples, they were vacuum-centrifuged until

394    approximately 3 μL of sample remained. Next, samples were resuspended in 6 μL of 5% ACN 0.1%

395    formic acid (FA).

396        Liquid chromatography coupled with tandem mass spectrometry (LC-MS/MS) was used to

397    analyse the samples, based on previously published protocols[9,50]. Samples were separated on a 15 cm

398    column (75 μm inner diameter in-house laser pulled and packed with 1.9 μm C18 beads (Dr. Maisch,

399    Germany) on an EASY-nLC 1200 (Proxeon, Odense, Denmark) connected to an Exploris 480

400    (CMNF-59632) or a Q-Exactive HF-X (Fontana Ranuccio *Stephanorhinus*) mass spectrometer (both

401    from Thermo Scientific, Bremen, Germany), with an integrated column oven. 4 (CMNF-59632) or 5

402    (FR sd - 295 - *Stephanorhinus*) μL of sample were injected. 0.1% FA in milliQ water was used as

403    buffer A and the peptides were separated with increasing buffer B (80% ACN 0.1% FA) with a 77

404    min gradient, increasing from 5% to 30% in 50 min, 30% to 45% in 10 min, 45% to 80% in 2 min,

405    and maintained at 80% for 5 min before decreasing to 5% in 5 min, and finally held for 5 min at 5%.

406    Flow rate was 250 nL/min. An integrated column oven was used to maintain the temperature at 40°C.

407        The two mass spectrometers were run with the same parameters except where specified, due

408    to changes in running software. Spray voltage was set to 2 kV, the S-lens RF level was set to 40%,

409    and the heated capillary was set to 275°C. Full scan mass spectra (MS1) were recorded at a resolution

410    of 120,000 at m/z 200 over the m/z range 350-1400. The AGC target value was set to 300% / 3e6

411    (Exploris / HF-X) with a maximum injection time of 25 ms. HCD-generated product ions (MS2) were

412    recorded in data-dependent top-10 mode and recorded at a resolution of 60,000. The maximum ion

413    injection time was 118 / 108 ms (Exploris / HF-X), with an AGC target value of 200% / 2e5 (Exploris

414    / HF-X). Normalised collision energy was set at 30 / 28% (Exploris / HF-X). The isolation window

415    was set to 1.2 m/z with a dynamic exclusion of 20 s. A wash-blank, using 5% ACN, 0.5% TFA, was

416    run between each sample and laboratory blank to limit cross-contamination.

417    Database construction

418        The protein reference alignment used in Cappellini *et al*. (2019)[9] was used as a starting point

419    to construct a database for sequence reconstruction. Due to the vast evolutionary distance between

420    CMNF-59632 and any extant taxa (>20 million years), a broader database was constructed to identify

421    sequence variants that may be known in other mammals. To construct a broader database, we

422    searched Uniprot and NCBI for each enamel protein, specifying the taxonomic grouping of 'Theria' to

423    include all therian mammals. To supplement available sequences, additional sequences were manually

424    extracted from available genomes, following the methodology from[51].

425        To investigate the relationships at the base of Rhinocerotidae, protein sequences translated

426    from *Elasmotherium sibiricum* genomic data[25] were generated.  To obtain the corresponding amino

427    acid sequences, we firstly collapsed the Paired-End (PE) reads and masked the conflict bases as "N"

428 using adapterRemoval[52]. Then, we mapped the collapsed reads against the reference genome of the

429 white rhinoceros (GCF_000283155.1_CerSimSim1) using the BWA MEM function[53] with the shorter

430 split hits being abandoned. After that, we removed duplicates using an in-house Perl script following

431 Liu et al. (2021)[25]. Finally, we extracted the gene sequences according to their locations on the

432 reference genome.

433 The remaining steps generally follow that outlined by[9]. To translate relevant genes, we used

434 ANGSD[54] to create consensus sequences for only those BAM files representing chromosomes with

435 genes of interest. To reduce the effects of post-mortem aDNA damage, we trimmed the first and last

436 five nucleotides from each DNA fragment. We formatted each consensus sequence as a blast

437 nucleotide database. To recover translated protein sequences, we performed a tblastn alignment[55],

438 with the corresponding *Ceratotherium simum* sequences as queries. Finally, we used ProSplign to

439 recover the spliced alignments, and ultimately, the translated protein sequences[56].

440 Protein identification

441 Thermo .raw files generated by the mass spectrometer were searched with various softwares

442 using an iterative search strategy to interpret spectra, characterise PTMs, and ultimately, reconstruct

443 protein sequences. For comparison, .raw files from a mediaeval ovicaprine (control) and an early

444 Pleistocene *Stephanorhinus* generated by Cappellini et al. (2019)[9] were also analysed. Among

445 samples from the Fontana Ranuccio *Stephanorhinus*, only FR4 was analysed.

446 We primarily employed MaxQuant[57] for sequence reconstruction and other downstream

447 aspects of data analysis. We performed two initial runs: 1) a more focussed run using the database we

448 modified from Cappellini et al. (2019)[9], and 2) a broad run using the 'Theria'-wide database we

449 constructed from publicly-available sequences.

450 In all runs, an Andromeda threshold of 40 and a delta score of 0 were set for both unmodified

451 and modified peptides. Minimum and maximum peptide lengths were specified as 7 and 25,

452 respectively. The default peptide false discovery rate (FDR) was used (0.01), while protein FDR was

453    increased to 1 to show possible low-abundance proteins. "Unspecific" digestion was specified. No

454    fixed post-translational modifications were set. Several PTMs were set as variable modifications in

455    our initial runs: glutamine and asparagine deamidation, methionine and proline oxidation, N-terminal

456    pyroglutamic acid from glutamic and aspartic acids, phosphorylation of serine, threonine, and

457    tyrosine, and the conversion of arginine to ornithine.

458    Proteins included in the database of common contaminants provided by MaxQuant (for

459    example proteinaceous laboratory reagents and human skin keratins), as well as reverse sequences,

460    were manually removed and not further examined. In addition, proteins detected in the laboratory

461    blank were also treated as contaminants, and not considered further.

462    To discover new SAPs and peptide variants not included in our database, we used additional

463    search tools. Peaks v. 7.0 wass used to attempt *de novo* sequencing and homology search was

464    performed using the SPIDER algorithm[58–60]. The open search capabilities of openPFind[61] and

465    MSFragger[62] were also used. When possible, the same settings were selected as in the MaxQuant

466    runs.

467    With our iterative search strategy, we integrated possible sequence variants from the results of

468    our *de novo*, homology searches, and open searches into hypothetical sequences from closely-related

469    taxa, to produce artificial sequences. These artificial sequences were included in a subsequent

470    MaxQuant search, and only incorporated into reconstructed sequences if identified and validated

471    using MaxQuant.

472    Sequence reconstruction + filtering

473    Prior to sequence reconstruction, all non-redundant PSMs were filtered using three criteria to

474    reconstruct only those peptide sequences and amino acid residues that we can confidently assign.

475    Sequences were accepted at two levels, resulting in two different datasets: 1) a minimally-filtered

476    dataset, and 2) a strictly-filtered dataset. This filtering starts with using Basic Local Alignment Search

477    Tool (BLAST)[63] to determine if peptides match any contaminants, beyond those included in

478  MaxQuant by default, such as soil bacteria and fungi. Next, MS/MS spectra are manually inspected

479  for each PSM to examine ion series coverage. At this stage, peptide sequences are accepted for the

480  strictly-filtered dataset only if each amino acid residue is covered (e.g., at least y-, b-, or a- ion

481  designates the mass of that specific amino acid, plus any identified PTMs) by at least two spectra,

482  following the approach outlined by Coutu et al. (2020)[64]. Additionally, for both strict- and minimally-

483  filtered datasets, poorly-supported spectra are also removed at this stage, and proteins are only

484  submitted for phylogenetic analysis if they are covered by at least two non-overlapping peptides.

485  Finally, under the strict-filtering criteria, BLAST is used again on any trimmed sequences, to remove

486  any that match contaminants.

487  Protein damage analysis

488  Characterization of protein degradation and post-translational modifications roughly follows

489  that of Cappellini et al. (2019)[9]. In addition to the primary run, three additional runs were performed

490  on each rhinocerotid sample, alongside the mediaeval control, to assess several different post-

491  translational modifications: 1) oxidative degradation of tryptophan, and included kynurenine ($\Delta M =$

492  $+3.994915$), oxolactone ($\Delta M = +13.979265$), tryptophandione ($\Delta M = +29.974178$) as variable

493  modifications, 2) oxidative degradation of histidine (His), including oxohistidine ($\Delta M = +15.9949$),

494  dioxohistidine ($\Delta M = +31.990$), His to hydroxyglutamate ($\Delta M = +7.979$), and His -to aspartic acid

495  ($\Delta M = -22.032$) as variable modifications, and 3) aromatic, including oxidation (WY) ($\Delta M =$

496  $+15.9949$) and dioxidation (WY) ($\Delta M = +31.990$) as variable modifications. Deamidation (NQ) was

497  also included as a variable modification for each run. After removing potential contaminants and

498  reverse hits, we used spectral counting to assess the extent of each PTM. In the case of arginine,

499  histidine, and tryptophan, we summed the total number of each amino acid residues in each sample,

500  after filtering for reverse hits and potential contaminants, and the total number of each modified

501  amino acid residue in each sample (Figure 2C and Ext. Figure 1BC). Deamidation levels were

502  calculated following the approach described in Mackie *et al*. (2018)[50] (Figure 2C) and site specific

503  rates were observed using the DeamiDATE algorithm[65].

504     Intra-crystalline protein decomposition analysis

505        Chiral amino acid analysis was undertaken on CMNF-59632 to evaluate the overall extent of

506     amino acid degradation in the intra-crystalline fraction of the enamel, enabling comparison to

507     previously-analysed specimens[9], and samples that had been experimentally-heated between 60 and

508     80°C for up to 17520 hours and to samples heated to 200 - 500 °C for up to 25 min. Enamel chips

509     were drilled using a Dremel ® 4000 (4000-1/45) drill with a diamond wheel point (4.4 mm (7105) by

510     Dremel ®) to remove any dentine which could be identified under a microscope (ZEISS Stemi 305,

511     Axiocam 105 R2). Samples were processed following the methods of Dickinson et al. (2019)[66]. To

512     remove excess powders, enamel chips were washed in deionized water and ethanol (Analytical-grade)

513     before being powdered with an agate pestle and mortar. Powdered samples were weighed into a single

514     plastic microcentrifuge tube and bleached (NaOCl-12%, 50 μL mg$^{-1}$ of enamel) for 72 hours to

515     remove the inter-crystalline amino acids and any contamination. This bleached sample was washed

516     five times with deionized water, and then once with methanol (HPLC-grade), before being left to dry

517     overnight.

518        The dried bleached sample was then divided into four subsamples: two for replicate analysis

519     of the free amino acids (FAA) and two for replicate analysis for the total hydrolysable amino acids

520     (THAA). The THAA subsamples were dissolved in HCl (7 M, 20 μL mg$^{-1}$, Analytical grade) in a

521     sterile 2 mL glass vial (Wheaton), purged with $N_2$ to reduce oxidation and heated at 110°C for 24 h in

522     an oven (BINDER GmbH series). The acid was then removed by centrifugal evaporation (Christ

523     RVC2-25). Then, THAA and FAA fractions were subjected to a biphasic separation procedure[66,67] to

524     remove inorganic phosphate from the enamel samples. HCl was added to both FAA (1 M, 25

525     μL mg$^{-1}$) and THAA (1 M, 20 μL mg$^{-1}$) fractions in separate 0.5 mL plastic microcentrifuge tubes

526     (Eppendorf), and KOH (1 M, 28 μL mg$^{-1}$) was added into the acidified solutions, which then formed

527     mono-phasic cloudy suspensions. Samples were agitated and then samples were centrifuged (13,000

528     rpm for 10 min, Progen Scientific GenFuge 24D) to form a clear supernatant above a gel. The

529     supernatant was removed, and dried by vacuum centrifugation. The concentration of the intra-

530    crystalline amino acids, and their extent of racemisation (D/L value) were then quantified using RP-

531    HPLC (Agilent 1100 series HPLC fitted with HyperSil C18 base deactivated silica column (5 μm, 250

532    x 3 mm) and fluorescence detector) following a modified method of Kaufman & Manley (1998)[68].

533         For the RP-HPLC analysis, samples were rehydrated with an internal standard solution (L-

534    homo-arginine (0.01 mM), sodium azide (1.5 mM) and HCl (0.01 M)) and run alongside standards

535    and blanks. A tertiary mobile phase system (HPLC-grade acetonitrile:methanol:sodium buffer (21

536    mM sodium acetate trihydrate, sodium azide,1.3 μM EDTA, pH adjusted to $6.00 \pm 0.01$ with 10%

537    acetic acid and sodium hydroxide)) was used for analysis. D and L peaks of the following amino acids

538    were separated: aspartic acid and asparagine (Asx); glutamic acid and glutamine (Glx); serine (Ser),

539    alanine (Ala), valine (Val), phenylalanine (Phe), isoleucine (Ile), leucine (Leu), threonine (Thr),

540    arginine (Arg), tyrosine (Tyr) and glycine (Gly). During preparation, asparagine and glutamine

541    undergo rapid irreversible deamination to aspartic acid and glutamic acid respectively[69] and hence

542    they are reported together as Asx and Glx respectively.

543    Phylogenetic Analysis

544         A time-calibrated phylogenetic tree was inferred with the Bayesian phylogenetic software

545    RevBayes v.1.2.1[70]; https://revbayes.github.io/) under a constant-rate Fossilised Birth Death (FBD)

546    model[71,72]. The dataset consisted of enamel proteome data for 16 perissodactyl species (10 extant and 6

547    extinct), totalling 7 proteins and 3446 amino acids. Phylogenetic analyses were performed with both

548    the strict-filtered and minimally-filtered sequences for CMNF-59632, to observe any topological

549    differences between the two datasets and assess if filtering is warranted. As no major differences were

550    observed, only the results from the 'strictly-filtered' dataset are discussed. The proteome dataset was

551    partitioned by protein. A GTR + I (General Time Reversible + Invariant sites) amino acid substitution

552    model—where stationary frequencies of the 20 amino acids and exchangeability rates among amino

553    acids are free to vary and estimated from data—was applied to each partition. Preliminary unrooted

554    phylogenetic analyses performed on each protein showed evidence for within-protein $\boldsymbol{\Gamma}$(Gamma)-

555    distributed rate variation only for MMP20, hence $\boldsymbol{\Gamma}$-distributed rate variation was modelled only for the

556    MMP20 partition. A relaxed clock model with uncorrelated lognormal-distributed rates (UCLN) was

557    applied to allow rate variation across branches. The prior on the average clock rate was set as a

558    loguniform distribution (min=$10^{-8}$, max=$10^{-2}$ substitutions per lineage/million year). The prior on the

559    clock rate standard deviation was set as an exponential distribution with mean equal to 0.587405,

560    corresponding to one order of magnitude of clock rate variation among branches. The FBD tree model

561    allows for placement of extinct species in a phylogenetic tree while simultaneously estimating the rates

562    of speciation, extinction, and fossilisation (sampling of species in the past). The priors on speciation,

563    extinction, and fossilisation parameters were set as uniform distributions bounded between 0 and 10.

564    The sampling probability for extant species was fixed to 0.5882353 (10/17), corresponding to the

565    fraction of extant perissodactyl species included in the analysis, and assuming uniform sampling of

566    extant taxa. The three species of Equidae in the analysis (*Equus caballus*, *E. przewalskii*, *E. asinus*)

567    were constrained as outgroup to other perissodactyls (Tapiridae + Rhinocerotidae). Tip ages of fossil

568    taxa were given a uniform prior distribution ranging from minimum to maximum age of the deposit

569    where each fossil has been found. The prior on the origin age of the tree was set as a uniform distribution

570    with minimum = 54 Ma, corresponding to the oldest fossil that can be unequivocally assigned to crown

571    Perissodactyla (*Cambaylophus vastanensis* from the early Ypresian Cambay Shale[73]), and maximum =

572    100 Ma, corresponding to the beginning of the Late Cretaceous and a very lax upper boundary on the

573    origin of placental mammals[74]. Additional constraints on node ages based on the fossil record of

574    perissodactyls were set to improve precision of divergence age estimates. Each node calibration was set

575    up as a soft-bounded uniform distribution with normally distributed tails, with 2.5% of the distribution

576    younger than the minimum age (allowing for potential misattribution of the oldest fossil of a clade) and

577    2.5% of the distribution older than the maximum age. Monophyly was not enforced when setting up

578    these node calibrations. The following age constraints have been applied to five nodes:

579    1) Node = crown Perissodactyla; soft minimum = 54 Ma, with the same justification as the minimum

580    on the origin age prior; soft maximum = 66 Ma, corresponding to the Cretaceous–Palaeogene boundary,

581    before which no unambiguous crown placental fossils are known. 2) Node = Rhinocerotina (crown

582    rhinoceroses); soft minimum = 22.6 Ma, corresponding to the earliest putative appearance of a crown

583  rhinoceros in the fossil record (*Gaindatherium* cf. *browni* from the Aquitanian upper member of the

584  Chitarwata Formation[75,76]; soft maximum = 44 Ma, corresponding to the minimum age of

585  Rhinocerotidae as supported by fossil and phylogenetic evidence [25]. 3) Node = Diceroti (*Ceratotherium*

586  + *Diceros*); soft minimum = 5.3 Ma, corresponding to the minimum age of the oldest deposits yielding

587  *Diceros bicornis* fossils (Lothagam and Albertine[77,78]; soft maximum = 7.3 Ma, as in Liu et al. (2021)[25].

588  4) Node = *Rhinoceros unicornis* + *Rhinoceros sondaicus*; soft minimum = 1.9 Ma, corresponding to the

589  early Pleistocene appearance of *Rhinoceros unicornis* in the fossil record[79,80]; soft maximum = 5.3 Ma,

590  as in Liu et al. (2021)[25]. 5) Node = *Dicerorhinus* + *Stephanorhinus* + *Coelodonta*; soft minimum = 13

591  Ma, corresponding to middle Miocene remains of *Dicerorhinus* from the middle Siwaliks of

592  Pakistan[25,81]; soft maximum = 22.6 Ma, corresponding to the oldest crown rhinoceros fossil as in the

593  soft minimum of calibration 2.

594  The Markov chain Monte Carlo (MCMC) was set up as 4 independent runs, running for 50,000

595  iterations and sampling every 10, averaging between 262.2 and 279.2 moves per iterations.

596  Convergence between runs was checked by visually inspecting and calculating effective sample sizes

597  (ESSs) of parameter estimates on Tracer v.1.7.2[82]. A maximum *a posteriori* (MAP) tree was calculated

598  to summarise the posterior distribution of trees, with 20% burn-in. In the analysis of the minimally-

599  filtered dataset, one of the 4 runs was discarded from the MAP tree calculation, as it converged only in

600  the last 10% of the MCMC.

# Methods References

602  41. Omar, G. *et al.* Fission-track dating of haughton astrobleme and included biota, devon island,

603  Canada. *Science* **237**, 1603–1605 (1987).

604  42. Young, K. E. *et al.* Impact thermochronology and the age of Haughton impact structure,

605  Canada. *Geophys. Res. Lett.* **40**, 3836–3840 (2013).

606  43. Stephan, T. & Jessberger, E. K. Isotope systematics and shock-wave metamorphism: III. K-

607  Ar in experimentally and naturally shocked rocks; the Haughton impact structure, Canada.

608  *Geochim. Cosmochim. Acta* **56**, 1591–1605 (1992).

609     44.Sherlock, S. C. *et al.* Re-evaluating the age of the Haughton impact event. *Meteorit. Planet.*

610     *Sci.* **40**, 1777–1787 (2005).

611     45.Erickson, T. M. *et al.* Resolving the age of the Haughton impact structure using coupled

612     40Ar/39Ar and U-Pb geochronology. *Geochim. Cosmochim. Acta* **304**, 68–82 (2021).

613     46.Hickey, L. J., Johnson, K. R. & Dawson, M. R. The stratigraphy, sedimentology, and fossils

614     of the Haughton formation: A post-impact crater-fill, Devon island, N.w.t., Canada. *Meteoritics*

615     **23**, 221–231 (1988).

616     47.Von Koenigswald, W., Holbrook, L. T. & Rose, K. D. Diversity and Evolution of Hunter-

617     Schreger Band Configuration in Tooth Enamel of Perissodactyl Mammals. *acpp* **56**, 11–32

618     (2011).

619     48.Taurozzi, A. J. *et al.* Deep-time phylogenetic inference by paleoproteomic analysis of dental

620     enamel. *Nat. Protoc.* (2024) doi:10.1038/s41596-024-00975-3.

621     49.Rappsilber, J., Mann, M. & Ishihama, Y. Protocol for micro-purification, enrichment, pre-

622     fractionation and storage of peptides for proteomics using StageTips. *Nat. Protoc.* **2**, 1896–1906

623     (2007).

624     50.Mackie, M. *et al.* Palaeoproteomic Profiling of Conservation Layers on a 14th Century Italian

625     Wall Painting. *Angewandte Chemie* vol. 130 7491–7496 Preprint at

626     https://doi.org/10.1002/ange.201713020 (2018).

627     51.Rüther, P. L. *et al.* SPIN enables high throughput species identification of archaeological

628     bone by proteomics. *Nat. Commun.* **13**, 2458 (2022).

629     52.Lindgreen, S. AdapterRemoval: easy cleaning of next-generation sequencing reads. *BMC Res.*

630     *Notes* **5**, 337 (2012).

631     53.Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.

632     *arXiv [q-bio.GN]* (2013).

633     54.Korneliussen, T. S., Albrechtsen, A. & Nielsen, R. ANGSD: Analysis of Next Generation

634     Sequencing Data. *BMC Bioinformatics* **15**, 356 (2014).

635     55.Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database

636     search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).

637    56. Kuznetsov, A. & Bollin, C. J. NCBI Genome Workbench: Desktop Software for Comparative

638    Genomics, Visualization, and GenBank Data Submission. *Methods Mol. Biol.* **2231**, 261–295

639    (2021).

640    57. Tyanova, S., Temu, T. & Cox, J. The MaxQuant computational platform for mass

641    spectrometry-based shotgun proteomics. *Nat. Protoc.* **11**, 2301–2319 (2016).

642    58. Han, Y., Ma, B. & Zhang, K. SPIDER: software for protein identification from sequence tags

643    with de novo sequencing error. *Proc. IEEE Comput. Syst. Bioinform. Conf.* 206–215 (2004).

644    59. Yuen, D., Zhang, W., Zhang, Z., Lajoie, G. A. & Ma, B. PEAKS DB: de novo sequencing

645    assisted database search for sensitive and accurate peptide identification. *& cellular proteomics*

646    (2012).

647    60. Ma, B. & Johnson, R. De novo sequencing and homology searching. *Mol. Cell. Proteomics*

648    **11**, O111.014902 (2012).

649    61. Chi, H. *et al.* Open-pFind enables precise, comprehensive and rapid peptide identification in

650    shotgun proteomics. *bioRxiv* 285395 (2018) doi:10.1101/285395.

651    62. Kong, A. T., Leprevost, F. V., Avtonomov, D. M., Mellacheruvu, D. & Nesvizhskii, A. I.

652    MSFragger: ultrafast and comprehensive peptide identification in mass spectrometry–based

653    proteomics. *Nat. Methods* **14**, 513–520 (2017).

654    63. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment

655    search tool. *J. Mol. Biol.* **215**, 403–410 (1990).

656    64. Coutu, A. N. *et al.* Palaeoproteomics confirm earliest domesticated sheep in southern Africa

657    ca. 2000 BP. *Sci. Rep.* **11**, 6631 (2021).

658    65. Ramsøe, A. *et al.* DeamiDATE 1.0: Site-specific deamidation as a tool to assess authenticity

659    of members of ancient proteomes. *J. Archaeol. Sci.* **115**, 105080 (2020).

660    66. Dickinson, M. R., Lister, A. M. & Penkman, K. E. H. A new method for enamel amino acid

661    racemization dating: A closed system approach. *Quat. Geochronol.* **50**, 29–46 (2019).

662    67. Dickinson, M. R. Enamel amino acid racemisation dating and its application to building

663    Proboscidean geochronologies. (University of York, 2018).

664    68. Kaufman, D. S. & Manley, W. F. A new procedure for determining dl amino acid ratios in

665     fossils using reverse phase liquid chromatography. *Quat. Sci. Rev.* **17**, 987–1000 (1998).

666     69.Hill, R. L. Hydrolysis of proteins. *Adv. Protein Chem.* **20**, 37–107 (1965).

667     70.Höhna, S. *et al.* RevBayes: Bayesian Phylogenetic Inference Using Graphical Models and an

668     Interactive Model-Specification Language. *Syst. Biol.* **65**, 726–736 (2016).

669     71.Heath, T. A., Huelsenbeck, J. P. & Stadler, T. The fossilized birth–death process for coherent

670     calibration of divergence-time estimates. *Proceedings of the National Academy of Sciences* **111**,

671     E2957–E2966 (2014).

672     72.Zhang, C., Stadler, T., Klopfstein, S., Heath, T. A. & Ronquist, F. Total-Evidence Dating

673     under the Fossilized Birth-Death Process. *Syst. Biol.* **65**, 228–249 (2016).

674     73.Kapur, V. V. & Bajpai, S. Oldest South Asian tapiromorph (Perissodactyla, Mammalia) from

675     the Cambay Shale Formation, western India, with comments on its phylogenetic position and

676     biogeographic implications. *Palaeobotanist* **64**, 95–103 (2015).

677     74.Carlisle, E., Janis, C. M., Pisani, D., Donoghue, P. C. J. & Silvestro, D. A timescale for

678     placental mammal diversification based on Bayesian modeling of the fossil record. *Curr. Biol.*

679     **33**, 3073–3082.e3 (2023).

680     75.Métais, G. *et al.* Lithofacies, depositional environments, regional biostratigraphy and age of

681     the Chitarwata Formation in the Bugti Hills, Balochistan, Pakistan. *J. Asian Earth Sci.* **34**, 154–

682     167 (2009).

683     76.Antoine, P.-O. *et al.* Chapter 16. Mammalian Neogene Biostratigraphy of the Sulaiman

684     Province, Pakistan. in *Fossil Mammals of Asia* 400–422 (Columbia University Press, 2013).

685     77.Pickford, M., Senut, B. & Hadoto, D. Geology and palaeobiology of the Albertine Rift valley,

686     Uganda-Zaire. Volume I: geology. *Publication occasionnelle-Centre international pour la*

687     *formation et les échanges géologiques* (1993).

688     78.Brown, F. H. & McDougall, I. Geochronology of the Turkana depression of northern Kenya

689     and southern Ethiopia. *Evol. Anthropol.* **20**, 217–227 (2011).

690     79.Tong, H. & Moigne, A.-M. Quaternary rhinoceros of China. *Acta anthropologica sinica* **19**,

691     257–263 (2000).

692     80.Antoine, P.-O. *et al.* A new rhinoceros clade from the Pleistocene of Asia sheds light on

693   mammal dispersals to the Philippines. *Zool. J. Linn. Soc.* **194**, 416–430 (2022).

694   81.Heissig, K. Palaeontologische und geologische untersuchungen im Tertiaer von Pakistan. V.

695   Rhinocerotidae (Mamm.) aus den unteren und mittleren Siwalik-Schichten. (1972).

696   82.Rambaut, A., Drummond, A. J., Xie, D., Baele, G. & Suchard, M. A. Posterior

697   Summarization in Bayesian Phylogenetics Using Tracer 1.7. *Syst. Biol.* **67**, 901–904 (2018).

698   83.Gessulat, S. *et al.* Prosit: proteome-wide prediction of peptide tandem mass spectra by deep

699   learning. *Nat. Methods* **16**, 509–518 (2019).

700

714   **Author contributions:** R.P., R.D.E.M, N.R., D.F., and E.C. designed the study. N.R, D.F. and M.G.

715   conducted fieldwork at the Haughton Crater site. R.D.E.M., N.R., D.F., M.G., R.F., and L.B.,

716   provided ancient samples. R.P., M.M, A.B., N.H., I.P., S.L., J.R-M., A.R., F.M., M.R.D., C.B., and

717   G.S. performed data generation and analysed data with support from E.C., S.H., E.W., N.R.,

718   R.D.E.M., and D.F.. R.P., M.M, and E.C. wrote the manuscript with contributions from all authors.

719  **Competing interests:** The authors declare no competing interests.

720  **Corresponding authors:** Correspondence to Ryan Sinclair Paterson, Danielle Fraser, or Enrico
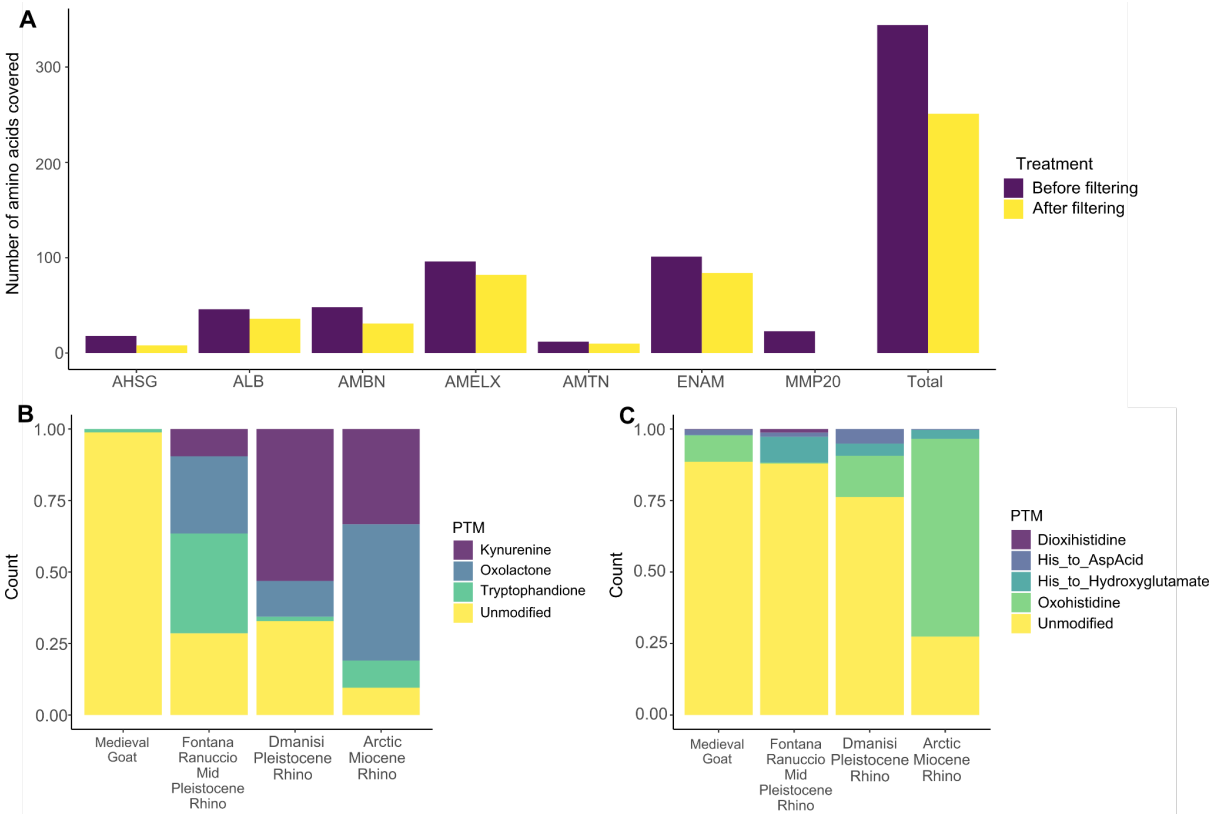
721  Cappellini.
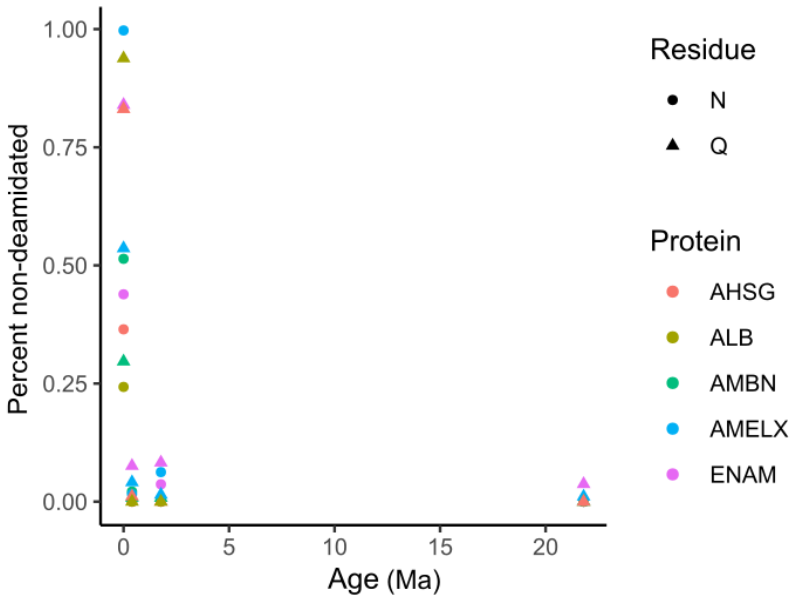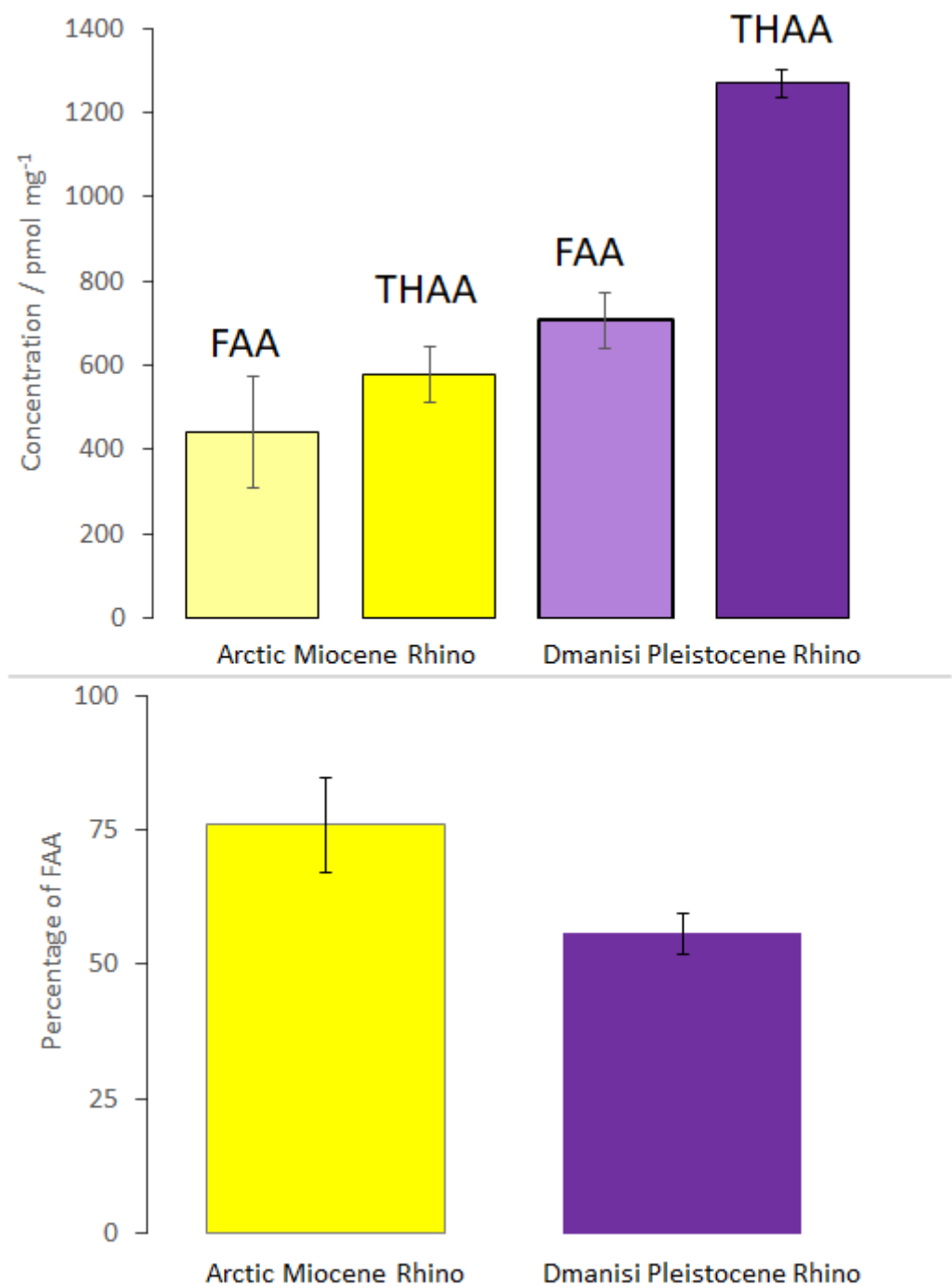
722

723

724

725

726

727

728

729

730

# Extended Data



**Extended Figure 1) Proteome preservation in CMNF-59632**. A) Amino acid count for each

identified protein, before and after filtering (see Methods), B) PTMs related to oxidative degradation

of tryptophan for CMNF-59632, compared to enamel proteomes from other ancient rhinos and a

mediaeval ovicaprine, C) PTMs related to oxidative degradation of histidine for same taxon set. The

moderate protein preservation in CMNF-59632, indicated by the lower amino acid coverage

compared to other ancient enamel proteomes, is further supported by the high incidence of PTMs

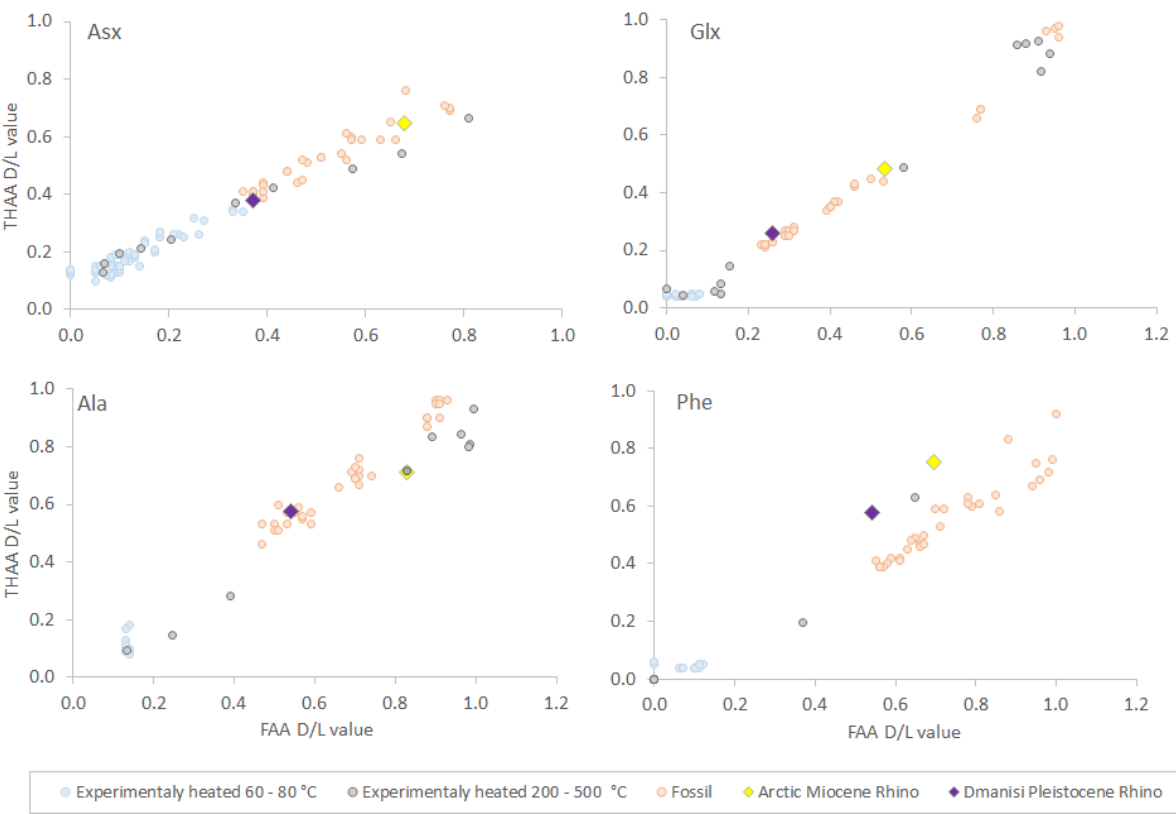related to oxidative degradation, compared to other fossil rhinocerotids.

740

**Extended Figure 2) Deamidation rates in fossil rhinocerotid enamel proteomes, plotted against geological age.** Data used is CMFN-59632 from Haughton Crater (21.8 Ma), DM.5/157 from Dmanisi (1.77 Ma), CGG 1_023342 from Fontana Ranuccio (0.4 Ma), and a mediaeval control sample (0.005 Ma). While useful for establishing authenticity of an ancient proteome, deamidation rates plateau relatively quickly, so they are not reliable for assessing relative degradative state in ancient proteomes from deep geological timescales.
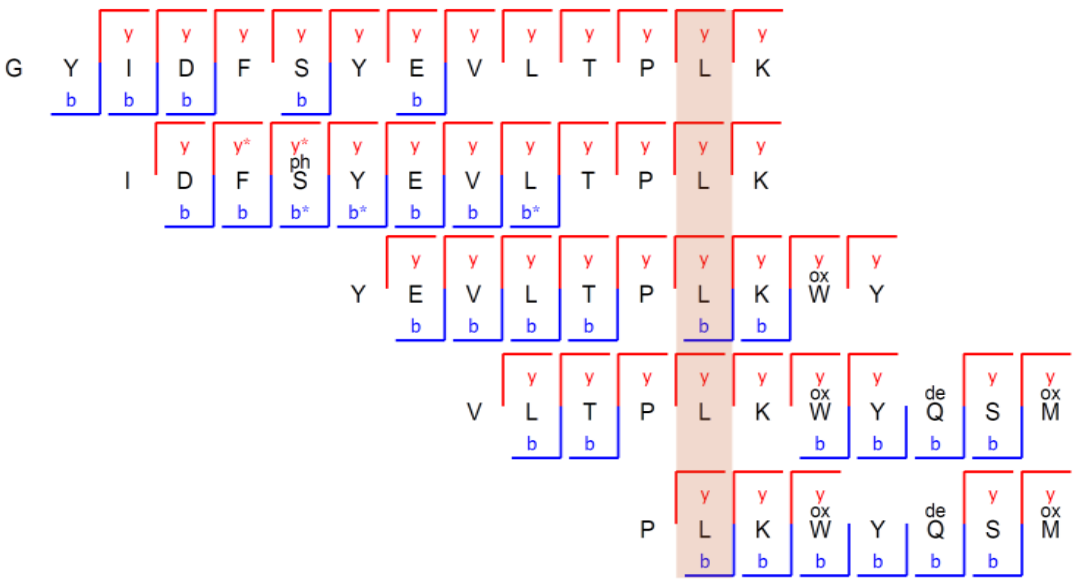
741

742

743

744

745

746

747

748

**Extended Figure 3) Comparison of FAA & THAA concentration (top) and %FAA (bottom) of the Arctic Miocene rhinocerotid with the Dmanisi Pleistocene *Stephanorhinus*[11].** Error bars represent 1 standard deviation about the mean for preparative replicates. The lower overall concentration, higher %FAA and yet incomplete hydrolysis in the Arctic Miocene rhino is consistent with endogenous peptides in the tooth enamel.

754

755



**Extended Figure 4) Asx, Glx, Ala and Phe FAA vs THAA D/L values for tooth enamel from Arctic Miocene rhino from Ellesmere Island Canada, and the Dmanisi Pleistocene rhino**. A data set consisting of published and unpublished enamel data from other rhino palaeontological and experimental data has been included for comparison. The good correlation between FAA & THAA for the Arctic Miocene rhino (CMNF-59632) sample supports the presence of closed system original peptides and their constituent amino acids in this Miocene sample.

763

764

765

766

**Extended Figure 5) Subset of overlapping peptides supporting a SAP at AMELX position 53.**
While all other later-diverging rhinocerotids, including *Elasmotherium*, display a valine (V) at
position 53 (highlighted in light red; following numbering of reference sequence A0A5F5PLN8),
CMNF-59632 displays a leucine (L) (or isoleucine (I)), representing the ancestral condition in
Perissodactyla. The peptide sequences depicted here represent just a small portion of the peptide-
spectrum matches covering this position. Together, these peptides display high Andromeda scores,
extended, in some cases complete, ion series, and the presence of several PTMs supporting their
endogeneity (phosphorylation) and ancientness (tryptophan oxidation, glutamine deamidation).