

1 **Convergent evolution through independent** 2 **rearrangements in the primate amylase locus**

3 Charikleia Karageorgiou,¹ Petar Pajic,² Stefan Ruhl,³ and Omer Gokcumen^{1,4*}

4 ¹Department of Biological Sciences, University at Buffalo, Buffalo, NY, USA.

5 ²Department of Chemistry, Yale University, New Haven, CT, 06511 USA.

6 ³Department of Oral Biology, School of Dental Medicine, University at Buffalo, Buffalo, NY, USA.

7 ⁴Lead contact

8 *Correspondence: omergokc@buffalo.edu

9 **SUMMARY**

10 Structurally complex genomic regions can foster evolutionary convergence by repeatedly generating
11 gene duplications that yield similar expression patterns and traits across lineages. Focusing on the
12 primate amylase locus, we leveraged high-quality genome assemblies from 53 primate species and
13 multi-tissue transcriptomes from Old World monkeys to reconstruct the evolutionary history of recurrent
14 duplications. We show that lineage-specific long terminal repeat retrotransposon insertions may be
15 associated with initial structural instability, while subsequent duplications are primarily driven by
16 non-allelic homologous recombination. Independent duplications in rhesus macaques, olive baboons, and
17 great apes produced distinct amylase copies with convergent expression in pancreas and salivary glands
18 and signals of episodic diversifying selection, consistent with emerging functional divergence. Our
19 analyses indicate that an ancestral gene with dual pancreas and salivary expression in Catarrhini
20 duplicated in great apes, facilitating subfunctionalization and regulatory rewiring. These findings illuminate
21 how modular structural and regulatory variation drives evolutionary innovation and molecular
22 convergence.

23 **KEYWORDS**

24 Amylase locus, Gene duplication, Recurrent duplication, Structural variation, Transposable elements,
25 Non-allelic homologous recombination (NAHR), Convergent evolution, Tissue-specific gene expression,
26 Subfunctionalization

27 **INTRODUCTION**

28 Convergent evolution, the independent emergence of similar traits in distinct lineages, has long intrigued
29 evolutionary biologists seeking to understand how similar phenotypes arise from different genetic starting
30 points^{1,2}. Recent advances in long-read sequencing have allowed us to characterize structurally complex
31 genomic regions accurately at the nucleotide resolution, revealing that these regions frequently undergo
32 recurrent structural variation³⁻⁶, including gene duplications, that can drive major shifts in biological
33 function⁷⁻¹¹. Integrating these insights, an emerging model posits that structurally complex loci, through
34 repeated gene duplications and regulatory rewiring, can serve as substrates for molecular convergence
35^{5,12,13}. In particular, such loci may produce similar spatial expression patterns of gene families in
36 phylogenetically distant lineages. Yet, the specific mutational mechanisms that give rise to these
37 duplications, the evolutionary forces shaping nucleotide variation among paralogs, and the processes by
38 which regulatory elements are reshuffled during structural rearrangements remain largely unresolved.
39 These questions are at the heart of understanding how structural and regulatory complexity contributes to
40 evolutionary innovation and the repeated emergence of similar traits across lineages.

41 The amylase locus is one of the most intriguing structurally complex regions in mammalian genomes,
42 notable for its exceptionally rapid structural evolution. It is one of the fastest-evolving loci in the human
43 genome, despite being essential in starch metabolism¹⁴. Amylase in mammals is primarily expressed in

44 the pancreas. However, in some lineages, the gene has undergone a regulatory shift to include
45 expression in the salivary glands^{15–17}. Variation in amylase copy number has been proposed to relate to
46 dietary starch intake in several lineages. In humans, salivary amylase (*AMY1*) copy number shows
47 population-level associations with starch-rich diets^{3,18,19}, albeit the evolutionary and functional
48 interpretation of this relationship remains an area of active investigation²⁰. Similar patterns of amylase
49 gene duplications and concordant expression shifts of the duplicated copies have evolved multiple times,
50 independently in different mammalian lineages, thus suggesting convergent mechanisms in response to
51 dietary shifts^{17,19,21,22}. However, the mutational mechanisms underlying the independent gene duplications
52 in the amylase locus and proximate regulatory sequences remain unexplored. Investigating these
53 processes, using this locus as a model, could provide key insights into fundamental aspects of genomic
54 evolution, including neofunctionalization, subfunctionalization, and gene expression dosage regulation.

55 In humans, saliva amylase is the most abundantly secreted enzyme in the oral cavity²³, with expression
56 levels 6–8-fold higher than in other great apes^{18,24,25}. This heightened expression was considered a
57 human-specific trait^{15,17,18}. However, across the primate phylogeny, other species also exhibit high salivary
58 amylase expression, including Old World monkeys²⁶ and capuchins¹⁷. The prevailing model suggests that
59 the ancestral amylase gene was expressed initially in the pancreas, and was duplicated independently in
60 different primate lineages, where some duplications acquired expression in the salivary glands. In apes,
61 where this process has been best studied, the shift to salivary expression of one *AMY1* duplicate has
62 been attributed to the insertion of an endogenous retroviral element (ERV) upstream of *AMY1*^{17,27–29}. In
63 humans, additional *AMY1* duplications led to increased saliva expression levels³⁰. In other nonhuman
64 primates, the genetic and regulatory mechanisms underlying salivary-gland-specific amylase expression
65 remain unknown.

66 Therefore, we compared the evolutionary history of the amylase locus in great apes and Old World
67 monkeys, offering an ideal framework for investigating how rapidly evolving, structurally complex regions
68 can give rise to convergent gene expression patterns. Specifically, we asked whether the recurrent gain
69 of salivary gland expression of *AMY* in these lineages is driven by shared mutational mechanisms and
70 regulatory shifts, or by distinct molecular events that lead to similar outcomes. To address this, we
71 leverage 244 high-quality primate genome assemblies and transcriptome data from pancreas, liver, and
72 salivary glands of rhesus macaques (*Macaca mulatta*) and olive baboons (*Papio anubis*). Using this data,
73 we identified features of the amylase locus associated with salivary gland expression, characterize the
74 mutational processes underlying lineage-specific gene duplications, and examine how structural and
75 regulatory variation is linked to *AMY* function and expression. In doing so, our study not only advances
76 understanding of amylase locus evolution and regulation in primates but also provides a broader model
77 for investigating how gene duplications in complex genomic regions can drive convergent evolution of
78 tissue-specific expression.

79 RESULTS

80 ***Ancestral and recurrent independent duplications shape primate amylase genetic structural*** 81 ***diversity***

82 Previous studies have documented extensive structural variation within the human amylase locus,
83 identifying multiple independent copy number changes and inversions^{3,4}. To determine if similar structural
84 complexity extends across non-human primates, we analyzed 244 primate genomes, successfully
85 curating continuous contigs spanning the amylase locus in 53 species (total 69 genomes, **Table S1 & 2**;
86 see **Methods**). For 16 species, multiple high-quality genome assemblies (ranging from 2 to 4 genomes
87 per species) exist, allowing us to assess within-species variation (**Table S3**). In the remaining 174
88 genomes, the amylase locus could not be completely resolved within a single contig, underscoring the
89 challenges of assembling this complex region, even with high-quality long-read data. Nevertheless, the
90 69 curated genomes provided a robust framework for reconstructing *AMY* structural variation across
91 primates, documenting the expansions and contractions of the copy number variation of *AMY* genes, and
92 predicting the ancestral structural states (**Figure 1A**, **Table S4**, see **Methods**).

93 Our results confirmed previous work^{17,27} that identified an ancestral duplication of the amylase locus in the
94 Catarrhini lineage (Catarrhini: the parvorder comprising Old World monkeys and apes; indicated by a
95 purple star in **Figure 1A**), lineage-specific duplications in New World monkeys, and a loss of a single
96 *AMY* gene in leaf-eating monkeys¹⁷. The recently available high-quality genome assemblies analyzed in
97 this study provide sequence-level resolution of lineage-specific copy number variations and led to the
98 discovery of novel duplication events. These include a burst of duplications in the orangutan lineage, a
99 previously noted¹⁸ but uncharacterized copy number increase in bonobos relative to chimpanzees
100 (**Figure 1B**), and recurrent independent duplications within Old World monkey genera, as described
101 below. In addition, we found copy number variation among lesser ape species (i.e., in gibbons) (**Figure**
102 **1A; Figure S1; Table S5**), yet it remains unclear whether this variation reflects an ancestral loss followed
103 by lineage-specific gains, independent lineage-specific losses, or incomplete lineage sorting (**Figure S2**),
104 posing an intriguing question for future research

105 Lemurs provide an ideal outgroup for studying the rest of the primate phylogeny. We found that all (n=11)
106 but one of the analyzed lemur species harbor a single *AMY* gene per haploid genome, a structure shared
107 with several non-lemur primates, suggesting that this configuration likely represents the ancestral
108 amylase haplotype in primates. One exception is *Microcebus murinus*, which carries an additional *AMY*
109 copy (**Figure 1C**). Taken together, our findings illustrate a complex evolutionary history of the amylase
110 locus characterized by recurrent independent duplications and losses, thus contributing a remarkable
111 structural diversity in primates. These patterns are consistent with the amylase region behaving as a
112 mutational hotspot across primates as defined by^{31–33}. More genomes from additional species and more
113 comprehensive dietary data would be required to robustly test the relationship between diet and *AMY*
114 copy number in primates. Our study sets the stage to more comprehensively analyze the evolution of
115 functional variation within a hotspot of genomic variation within closely related primate species with
116 distinct dietary habits.

117 **Reconstruction of the amylase duplications in catarrhini**

118 To further elucidate the mutational basis of *AMY* structural variation in primates, we closely examined the
119 rhesus macaque and olive baboon genomes, contrasting their evolutionary histories with recent findings
120 from human haplotypes⁴. Previous studies leveraged Old World monkeys as an outgroup to explore *AMY*
121 copy number expansions in the human lineage²⁹. These analyses successfully identified an ancestral
122 Catarrhini duplication event, followed by a great ape-specific duplication featuring a lineage-specific 5'
123 ERV retrotransposon associated with salivary expression²⁷ (**Figure 2**). In addition, recent assemblies
124 available for these species suggested lineage-specific duplications within rhesus macaque and olive
125 baboon genomes that were undetected in earlier work.

126 Our findings provide phylogenetic evidence indicating that the primate ancestor possessed a single
127 amylase gene, orthologous to human *AMY2B* (**Figure 1A**). In the Catarrhini ancestor, after the
128 divergence of new world monkeys, an insertion of the 3' untranslated region of a γ -actin pseudogene
129 insertion occurred 5' upstream of the ancestral amylase gene, followed by the duplication of the
130 γ -actin-*AMY2B* segment, thereby generating a new amylase copy (*AMY1'*) (**Figure 2**). Later in the great
131 ape lineage, an endogenous retrovirus (ERV) was inserted into the γ -actin region flanking *AMY1'*. This
132 combined γ -actin-ERV-*AMY1'* segment duplicated into *AMY1* and the precursor of *AMY2A* in great
133 apes. Eventually, the progenitor of *AMY2A* underwent an ectopic deletion of a portion of the ERV
134 element, leading to its current structure as described further along with mutational mechanisms (**Figure**
135 **2**). Our observations within the great ape lineage builds upon previous findings reported by Meisler &
136 Ting²⁹ but provides a more granular picture. The observation that *AMY1'* represents the ancestral copy
137 leading to great ape *AMY2A* and *AMY1* is remarkable because these genes have distinct functions with
138 specialized expression in pancreas and salivary glands, respectively.

139 Lineage-specific duplications in rhesus macaques:

140 In Old World monkeys, we found that baboons and rhesus macaques retain both the ancestral *AMY2B*
141 gene as well as the derived *AMY1'* gene, which originated in the Catarrhini ancestor (**Figure 2A**).
142 Additionally, we identified one lineage-specific duplicate in macaques and two additional, species-specific

143 duplicates in baboons (**Figure 1A**). To better understand these events, we assessed the gene orthologies
144 of these duplicates and resolved the duplication breakpoints relative to the ancestral Catarrhini and great
145 ape haplotypes (**Figure S3**).

146 The novel lineage-specific duplication in rhesus macaques, which we termed *AMY_m*, is located between
147 *AMY2B* and *AMY1'* (**Figure 3A**). Sequence comparisons indicated that *AMY_m* exhibits the highest
148 similarity to *AMY2B*, suggesting its origin from the ancestral *AMY2B* gene. *AMY1'* in rhesus macaques on
149 the other hand, shares orthology with great apes' *AMY2A* and *AMY1*, although without clear one-to-one
150 orthology with either gene, suggesting that *AMY1'* represents the ancestral state from which *AMY1* and
151 *AMY2A* evolved.

152 By comparing the ancestral Catarrhini haplotype (*i.e.*, two-copy haplotype in **Figure 3A**) (see **Methods**)
153 with the rhesus macaque (*Macaca mulatta*) haplotype, we resolved the breakpoints of the duplicated
154 sequences (**Figure 3B**). We identified non-allelic homologous recombination (NAHR) as the primary
155 mechanism driving the segmental duplications. NAHR, in this case, is characterized by two key
156 signatures: first, *AMY_m* is flanked by segmental duplications, *AMY2B* at the 5' end and *AMY1'* at the 3'
157 end; and second, the duplicated segment containing *AMY_m* exhibits mosaic characteristics derived from
158 both flanking regions (**Figure S3**). Additionally, the duplication breakpoints overlap with the γ -actin
159 insertion element, suggesting a possible role of this element in having facilitated the duplication events
160 via NAHR. Among macaques, we were further able to determine that the duplication is specific to *sinica*
161 and *fascicularis* groups, and must have occurred after the divergence of the *silenus* group (**Figure 3B**),
162 dating this duplication to the relatively tight window of 4.5 to 5 million years ago based on the previously
163 published phylogenetic dating of this clade³⁴. Based on these findings, we constructed what we believe to
164 be the most plausible model that explains the mutational mechanism underlying this macaque-specific
165 duplication event (**Figure S4**).

166 Species-specific duplications in olive baboons:

167 We conducted a similar comparative analysis of the olive baboon amylase locus to identify shared and
168 lineage-specific amylase gene duplications and the mechanisms through which they arise. Specifically,
169 we investigated the structural configuration of the amylase locus by comparing the olive baboon
170 haplotype carrying four gene copies to the ancestral Catarrhini haplotype. We identified two
171 lineage-specific amylase copies in *Papio anubis* flanked by *AMY2B* on the 5' and *AMY1'* on the 3'
172 (**Figure 3C**), which we termed *AMY_{p1}* and *AMY_{p2}*.

173 These two novel genes were absent from other Old World monkey genomes, consistent with independent
174 duplication events specific to the *Papio* lineage. To infer the sequential order of these duplications, we
175 first decomposed the olive baboon amylase locus into four duplicated segments (see **Methods**). These
176 analyses showed that the segment harboring *AMY_{p1}* is a mosaic of the ancestral *AMY2B*- and
177 *AMY1'*-containing blocks, whereas the segment carrying *AMY_{p2}* is nearly identical to the *AMY_{p1}*
178 segment, consistent with a second NAHR event duplicating a pre-existing *AMY_{p1}*-containing block. An
179 initial inspection of the Guinea baboon (*Papio papio*) amylase locus revealed six identical tandem
180 amylase segments in addition to the ancestral *AMY2B* and *AMY1'* (see **Methods** and **Figure S5**), thus
181 we treated these redundant copies as unreliable for ordering the duplication events. Instead, we relied on
182 the internal segmental architecture of the olive baboon locus, together with competitive mapping among
183 the four olive baboon *AMY* paralogs, to establish that *AMY_{p1}* represents the first duplication derived from
184 the ancestral *AMY2B*-*AMY1'* configuration, whereas *AMY_{p2}* arose later as a second, likewise olive
185 baboon-specific duplication, both highlighted in orange in **Figure 3C**.

186 Based on the data sources available to us and our findings outlined above, the most parsimonious
187 mutational model for these olive baboon-specific duplications to be as follows: The ancestral haplotype,
188 containing *AMY2B* and *AMY1'*, underwent an initial non-allelic crossover between these two genes, giving
189 rise to the novel *AMY_{p1}* copy. Subsequently, haplotypes carrying *AMY2B*, *AMY_{p1}*, and *AMY1'*
190 experienced a second recombination event, this time between *AMY2B* and *AMY_{p1}*, resulting in the
191 emergence of the second novel copy, *AMY_{p2}*, which is now positioned between *AMY2B* and *AMY_{p1}* in
192 the extant olive baboon haplotype (**Figure 3C**). Collectively, our reconstructions indicate that *AMY_{p1}*

193 predates *AMYp2* and that both duplications arose after the complete split from Guinea baboons (*Papio*
194 *papio*). The divergence between olive and Guinea baboons has been dated to approximately 1.85 million
195 years ago based on previously published estimates^{35,36}, indicating that both *AMYp1* and *AMYp2* are
196 younger than this split and originated within the olive baboon lineage.

197 Taken together, our findings highlight lineage-specific duplications in rhesus macaques and olive baboons
198 that occurred independently of one another through NAHR and within approximately 10 million years of
199 evolutionary divergence. Two major questions remain, which we address in the next two sections. First,
200 what was the initial mutational driver of these duplications? Second, what could be the adaptive and
201 functional impact of these lineage-specific duplications?

202 Mechanistic inference of NAHR

203 Our inference of NAHR is based on multiple convergent lines of evidence assessed against the
204 established diagnostic criteria for NAHR, MMBIR and NHEJ^{37–39}. First, each lineage-specific duplicate
205 (*AMYm* in macaques; *AMYp1* and *AMYp2* in baboons) is a chimeric copy whose 5' and 3' blocks derive
206 from different flanking paralogues, with a single transition point, the hallmark of a strand exchange
207 between misaligned paralogous sequences (**Figure S3**). Second, pairwise alignments reveal that the
208 flanking homologous tracts far exceed the ~300-500 bp minimal efficient processing segment (MEPS)
209 established for meiotic NAHR⁴⁰, providing ample substrate for homologous recombination. Third, the
210 recurrence of independent events at the same flanking segments across macaques, baboons, bonobos
211 and great apes is itself diagnostic; NAHR characteristically produces clustered, recurrent breakpoints
212 mediated by a fixed pair of segmental duplications, whereas replication-based mechanisms (MMBIR)
213 generate non-recurrent rearrangements with unique breakpoint footprints^{37,39}. Finally, we observe neither
214 the complex junction architecture, embedded triplications, inversions, or 2-5 bp microhomology-supported
215 junctions, expected under MMBIR nor the short insertions, short duplications at the break junctions or
216 blunt joins without flanking homology that are characteristic of NHEJ³⁸. These observations collectively
217 and robustly support NAHR as the main driver of structural variation in the primate amylase locus.

218 ***The abundance of long terminal repeats (LTRs) correlates with gene copy number gains in*** 219 ***the amylase locus***

220 The amylase locus is exceptional in that it has evolved rapidly through independent structural
221 rearrangements across the tree of life, including in species such as fruit flies, mice, rats, and dogs^{17,22,41,42}.
222 Several studies have linked this extensive variation to dietary adaptations^{17,18,22}. As described above, the
223 mutational mechanisms driving additional duplications, following the initial structural changes in the
224 Catarrhini ancestors, have been characterized as non-allelic homologous recombination (NAHR).
225 However, the origins of primary duplications, which arose from a single-copy ancestral haplotype (likely
226 representing the ancestral state in primates), remain poorly understood.

227 Our previous work¹⁷ identified lineage-specific insertions of transposable elements coinciding with *AMY*
228 gene duplications across various mammalian taxa. The association between transposable elements
229 (TEs) and structural variation has been explored in multiple contexts, and recent studies have shown that
230 TE-mediated rearrangements can arise through diverse molecular mechanisms and induce structural
231 instability^{43,44}. Building on these observations, we hypothesized that TE elements might contribute to the
232 formation of primary duplications in primates, thereby predisposing the amylase locus to structural
233 instability. To test this hypothesis, we used a standardized primate repeat dataset to annotate
234 transposable elements across 53 primate genomes, thereby minimizing genome-specific biases (**Table**
235 **S6**; see **Methods** for details).

236 Our analyses revealed a wide range of transposable element content in the amylase locus across
237 species, from 25.23% in the northern greater galago (*Otolemur garnettii*) to 59.52% in the Bornean
238 orangutan (*Pongo pygmaeus*) (**Figure 4A**). In contrast to our expectations, we found that the amylase
239 locus is generally depleted in transposable elements compared to genome-wide averages; nevertheless,
240 we observe an enrichment in LTRs (**Figure 4B & 4C**). The general depletion is primarily driven by

241 reduced representation of common, active, short retrotransposons such as Alu sequences (**Figure 4B**),
242 suggesting that the locus is not broadly permissive to transposable element retention.

243 While other common transposable elements were underrepresented within the amylase locus in most
244 primate species, we observed an enrichment of LTR transposons (**Figure 4**, panels **C** and **D**). The
245 proportion of LTR elements in the locus, even after accounting for phylogenetic relatedness, correlates
246 with the number of amylase genes, ($p < 10^{-5}$, $R^2 = 0.62$; **Figure S6**). This correlation is most significant
247 within Catarrhini, where the majority of amylase gene gains and losses were observed, supporting the
248 hypothesis that LTRs contribute to structural instability in the amylase locus. If indeed LTRs are driving
249 the structural changes, we would expect to find LTRs proximal to the breakpoint junctions, as exemplified
250 by the *AMYp2* duplication in olive baboons (**Figure S7**). We also detected a strong negative correlation
251 between amylase gene copy number and the presence of DNA transposons, particularly *Charlie* and
252 *Tigger* elements ($p = 1.4e^{-6}$ and $p = 0.0005$ respectively; **Figure S8**), further raising the possibility that
253 non-LTR transposons might have limited retention in the locus. Together, these findings point to a
254 transposable element-specific pattern of enrichment and depletion within the amylase locus, with these
255 dynamics showing strong correlations with amylase gene copy number variation across primates.

256 To further assess the relationship between LTRs and segmental duplications, we identified orthologous
257 LTR insertions across primates, clustered them into orthogroups, and reconstructed ancestral
258 presence/absence states on the primate phylogeny (**Figure S9 & S10**). This analysis revealed elevated
259 LTR gain events on the branch leading to Catarrhini, coinciding with the initial amylase duplication events
260 in this lineage. Importantly, for the LTRs discussed above, the structural evidence indicates that the LTR
261 copies were part of the ancestral segment that was subsequently duplicated by NAHR, rather than
262 secondary insertions into pre-existing duplicates. Although we cannot conclusively demonstrate LTRs as
263 the direct cause of structural rearrangements, we hypothesize that their insertion may contribute not only
264 to structural instability but may also facilitate functional changes in regulatory regions as suggested
265 earlier²⁹. Due to their role in regulatory rewiring, their retention may generate long stretches of
266 homologous sequence, which in turn may cause subsequent NAHR events. Such dynamics in primates
267 may explain both the recurrence of amylase gene duplications, as well as their functional changes in
268 tissue expression.

269 **Signals of selection suggest potential functional variation among primate amylase paralogs**

270 While human amylase paralogs in general appear to evolve under negative selection without significant
271 amino acid divergence from each other⁴, recent studies highlighted potentially functional variations
272 among human paralogs⁴⁵. To identify such functional differences and signals of potential positive
273 selection, we aligned the coding sequences of all identified paralogs from olive baboons and rhesus
274 macaques with paralogs from ape genomes.

275 Our analysis yielded three major insights (**Figure 5A & 5B**). First, we identified a significant positive
276 selection signal on the internal branch leading to baboon and rhesus macaque paralogs compared to
277 great ape paralogs (aBSREL, $p = 0.038$). Second, we detected a premature stop codon mutation in the
278 ancestral *AMY2B* gene copy within baboons, previously unannotated in the NCBI database (**Figure 5A**).
279 Third, we detected strong evidence for positive selection specifically on *AMYp2* in olive baboons,
280 supported by both aBSREL ($p < 10^{-6}$) and RELAX analyses ($p < 0.0001$). Additionally, we identified six
281 codons exhibiting episodic positive selection that likely contribute to this overall selection signal
282 (**Methods; Figure S11**). The concurrent presence of a premature stop codon in *AMY2B* and positive
283 selection in *AMYp2* is remarkable. It has been shown previously that following gene duplications, one of
284 the copies may become pseudogenized and the other one retains the original function^{46,47}. Thus, it is
285 plausible here that the newly derived gene (*AMYp2*) functionally compensates for the observed loss of
286 function in the ancestral gene (*AMY2B*).

287 Next, we analyzed the functional annotations of primate amylase amino acid sequences, including known
288 glycosylation sites, streptococci-binding motifs, catalytic and active sites, and calcium-binding domains
289 (**Figure 5C**, see **Methods**). To evaluate whether the observed amino acid substitutions might affect
290 protein structure or function, we generated AlphaFold2 models for the amylase gene products of each

291 paralog, with a particular focus on the baboon-specific *AMYp2* sequence, which showed the strongest
292 signal of positive selection (**Figure S12**). The functional predictions suggest that most of the positively
293 selected substitutions do not overtly disrupt the protein fold, glycosylation motifs, or key catalytic
294 residues. However, subtle impacts on substrate affinity or protein stability cannot be ruled out (**see Table**
295 **S7** for predicted pathogenic and neutral substitutions). Notably, we identified one strongly selected site (p
296 = 0.01, MEME) in *AMYp2*, involving a threonine-to-serine substitution at the position 178 predicted as an
297 active site by NLM's conserved domain database⁴⁸, which may reflect fine-tuned functional divergence in
298 the olive baboon lineage. Collectively, our findings suggest that, unlike reported for humans⁴, Old World
299 monkey amylase gene duplications involve significant amino acid diversification, indicative of potential
300 neofunctionalization.

301 **Reconstructing the evolution of salivary gland-specific expression**

302 To investigate the expression patterns of lineage-specific amylase gene paralogs, we generated
303 transcriptome data from parotid, sublingual, and submandibular salivary gland, as well as liver and
304 pancreas tissues from six rhesus macaques and five olive baboons (**Table S8**). We also leveraged our
305 previously published transcriptomic data from the corresponding human salivary gland tissues⁴⁹, and
306 added data for pancreas and liver tissues from the GTEx database. Taking advantage of sequence
307 differences among paralogs within each species, we used Kallisto⁵⁰ to quantify transcript abundance of
308 each gene, as it offers high sensitivity and accuracy for distinguishing among closely related gene
309 copies⁵⁰ (see **Methods** for more detailed discussion). These findings led us to form several hypotheses to
310 explain the regulation of amylase duplicates.

311 The first clear pattern we observed in both baboons and rhesus macaques is that, similar to humans, the
312 last gene (at the 3' end) in the amylase cluster consistently shows elevated expression (parotid vs.
313 pancreas \log_2FC : rhesus *AMY1'*=6.08, baboon *AMY1'*=1.91) in salivary tissues relative to the other
314 paralogs (**Figure 5D**). However, in contrast to humans, where *AMY1* is exclusively expressed in saliva,
315 this gene in baboons and rhesus macaques retains expression (pancreas vs. liver \log_2FC : rhesus
316 *AMY1'*=0.40, baboon *AMY1'*=3.48) in the pancreas and liver, indicating broad expression patterns
317 (**Figure 5E**). Additionally, the relative contribution of each amylase paralog to total salivary expression
318 differs among species (**Figure 5D**). In humans, *AMY1* accounts for nearly all salivary gland expression
319 (parotid vs. pancreas \log_2FC : *AMY1A*=14.71, *AMY1B*=9.74, *AMY1C*=14.35). The lineage-specific
320 duplications *AMYm* and *AMYp2* also contribute to the overall salivary expression in rhesus macaques
321 and olive baboons, respectively.

322 Given that *AMY1'* in Old World monkeys represents the ancestral gene from which the great ape *AMY1*
323 and *AMY2A* genes evolved (**Figure 2A**), it offers a valuable framework for investigating how
324 tissue-specific expression arose in these newly duplicated genes. In humans, *AMY1* is expressed
325 exclusively in salivary glands, while *AMY2A* is expressed only in the pancreas (**Figure 5E**). Our
326 transcriptomic analyses in rhesus macaques and baboons shows that *AMY1'* is expressed in both
327 pancreas and salivary glands, likely reflecting the ancestral state for Catarrhini (**Figure 5D & 5E**). Based
328 on these observations, the most parsimonious explanation is that the ancestral *AMY1'* gene had already
329 acquired expression in both the pancreas and salivary glands, particularly the parotid gland, prior to the
330 divergence of great apes. Following duplication in the great ape lineage, subfunctionalization occurred:
331 *AMY1* retained salivary gland expression, while *AMY2A* lost salivary expression and became restricted to
332 the pancreas. This shift may have been facilitated by an ERV insertion, as previously proposed^{27,29}.
333 Together, these findings support a subfunctionalization model for the great ape *AMY1* and *AMY2A*
334 following their divergence from *AMY1'*.

335 **Multifactorial regulation of tissue-specific expression of amylase paralogs**

336 To investigate how regulatory elements contributed to the evolution of amylase gene expression in
337 primates, we examined transcription factor binding sites across paralogs and species. Using *in silico*
338 predictions of transcription factor binding sites and promoter regions for amylase paralogs in rhesus
339 macaques, olive baboons, and humans, we identified 262 distinct binding sites (consisting of 108 unique
340 TFBS) across the promoters of the nine annotated amylase gene paralogs analyzed (**Table S9**).

341 A simplistic model of tissue-specific expression assumes that the presence of a specific transcription
342 factor binding site determines tissue specificity. If this were the case, we would expect the transcription
343 factor binding motifs among primate amylase paralogs to cluster according to their expression in either
344 salivary glands or pancreas. However, we found no such pattern: paralogs with known tissue-specific
345 expression did not show consistent enrichment of salivary- or pancreas-biased motifs (**Figure 6**). Instead,
346 our results suggest a partial conservation in promoter binding site composition among primate paralogs,
347 where all three rhesus macaque paralogs (*AMY2B*, *AMYm*, and *AMY1'*) share a similar binding site
348 profile with olive baboon *AMY1'* and *AMYp2*. In contrast, the promoter regions of the pseudogenized olive
349 baboon *AMY2B* and the human paralogs (*AMY2B*, *AMY2A*, and *AMY1s*) exhibit a distinct transcription
350 factor binding sequence composition, consistent with the lineage-specific evolution of these gene copies
351 (**Figure 2A**).

352 FOXC1 is a key transcription factor involved in salivary gland development and expression⁵¹. Notably, we
353 found that all primate *AMY* paralogs, regardless of whether they are salivary gland- or pancreas-biased,
354 contain FOXC1 binding sites, with the exception of human *AMY2B* and the pseudogenized olive baboon
355 *AMY2B*. This observation suggests that the presence of a salivary-biased regulatory motif alone is
356 insufficient to dictate tissue-specific expression. Instead, expression patterns of *AMY* genes may be
357 further shaped by context-dependent factors such as chromatin accessibility, competitive binding, or the
358 presence of co-regulators. For example, despite containing a FOXC1 binding site, human *AMY2A* is
359 primarily expressed in the pancreas, highlighting the complexity of regulatory control at this locus.

360 These findings suggest that the presence or absence of salivary gland or pancreas-biased transcription
361 factor binding sites alone does not fully account for the observed tissue-specific expression patterns,
362 highlighting the potentially combinatorial nature of transcriptional regulation, an observation consistent
363 with previous findings from other tissue-specific systems⁴⁹.

364 DISCUSSION

365 In this study, we used the amylase gene locus as a model to investigate within the primate lineage how
366 structurally complex loci contribute to molecular convergence. Convergent evolution, where similar traits
367 arise independently across lineages, is a hallmark of adaptive evolution and highlights non-random
368 patterns shaped by natural selection. Structurally complex regions have recently emerged as key players
369 in this process, and it is thought that independent gene duplications might play an important role in this
370 context. Focusing on the amylase locus, which is one of the most structurally dynamic regions in the
371 human genome⁴, we expanded prior work in humans to include the broader primate phylogeny. We
372 identified multiple lineage-specific amylase gene duplications, including previously uncharacterized
373 expansions in bonobos, orangutans, lemurs, and New World monkeys. This comprehensive analysis of
374 amylase evolution enabled us to dissect the interplay between mutational mechanisms, positive selection,
375 and regulatory divergence underlying functional convergence.

376 Our findings revealed that NAHR-mediated rearrangements in baboons, macaques, and humans
377 occurred through distinct breakpoints, underscoring the recurrent nature of duplication events. Given that
378 NAHR relies on existing homologous sequences, we asked how primary duplications arise in the first
379 place. We observed a strong correlation between amylase copy number and LTR abundance across
380 species, consistent with the hypothesis that lineage-specific transposable element insertions may
381 contribute to the homology required for initiating structural instability. While we cannot exclude a
382 bidirectional relationship, in which segmental duplications create a permissive genomic context for
383 additional TE accumulation, the synteny-based evidence is consistent with LTRs providing additional
384 stretches of shared sequence that facilitate NAHR. This provides preliminary evidence for a broader
385 hypothesis, that transposable elements may be a plausible contributor in priming structurally complex loci
386 for recurrent evolution, a compelling direction for future studies across mammalian genomes.

387 Beyond structural variation, we detected signatures of positive selection at specific codons, suggesting
388 functional divergence among amylase paralogs in a species-specific manner. Such changes may be of
389 functional relevance because they can alter the conformation of the protein and change or create motifs
390 governing posttranslational modifications. Evidence supports that human *AMY* paralogs differ in

391 glycosylation potential and may influence oral microbiota composition⁵². The interplay of nucleotide-level
392 and structural variation lays the groundwork for future population-level studies in closely related species
393 with differing diets or pathogen pressures, offering a powerful framework to investigate the adaptive
394 relevance of *AMY* variation.

395 By integrating genomic variation, gene expression, and transcription factor motif analyses, we show that
396 *AMY* paralogs have independently evolved tissue-specific expression, particularly in salivary glands.
397 Rather than relying on distinct, tissue-specific transcription factors, expression biases appear to result
398 from varying combinations of binding motifs among paralogs and across species. This regulatory rewiring
399 is linked to structural rearrangements at the locus, including duplications, inversions, and deletions, that
400 reshape the genomic context of regulatory elements. Our findings, thus, highlight the likely importance of
401 distal regulatory elements beyond core promoters in mediating tissue-specific expression. Future
402 functional assays such as ATAC-seq, ChIP-seq, or Fiber-seq will be essential to fully resolve the evolving
403 regulatory landscape of the amylase locus.

404 A vivid example of the co-evolution of gene duplication and regulatory rewiring is seen in the great ape
405 *AMY2A* and *AMY1* genes. Our results show that these paralogs arose from an ancestral gene (*AMY1'*)
406 with dual expression in the pancreas and salivary glands. Following duplication, *AMY2A* and *AMY1*
407 subfunctionalized into largely pancreas- and salivary gland-specific expression, respectively. This
408 instance of regulatory and functional co-divergence exemplifies the distinct evolutionary innovation at the
409 primate amylase locus, ranging from fixed duplications (as in macaques and baboons), to extensive
410 intraspecific variation (as in humans), and even gene loss (in leaf-eating monkeys).

411 Notably, we found no species lacking the amylase gene entirely, suggesting that while the locus tolerates
412 considerable structural and regulatory flexibility, it remains under consistent functional constraint. This
413 balance, between mutational plasticity and adaptive necessity, places the amylase locus in what Ponting⁵³
414 has termed the “evolutionary twilight zone”. It is within this zone that convergence has repeatedly
415 emerged during primate evolution, driven by distinct molecular mechanisms acting on a structurally
416 complex genomic landscape.

417 We propose that structurally complex loci across mammalian genomes are hotspots of molecular
418 convergence, harboring exceptional evolutionary potential. Their structural complexity not only promotes
419 gene copy number expansion and contraction but also enables regulatory innovation, fueling expression
420 divergence and functional diversification of duplicated genes.

421 **Limitations of the study**

422 Several limitations of this study should be noted. First, our analyses rely on a limited number of genome
423 assemblies (at most 4) per species. While parsimony-based reconstruction of copy number evolution
424 shows clear phylogenetic concordance and we confirmed structural configurations on multiple assemblies
425 for 16 species, we cannot exclude within-species copy number polymorphism in the amylase locus.
426 Population-level surveys, analogous to what has been achieved for humans^{3,4,20}, will be essential to
427 characterize the full extent of structural variation at this locus across primates.

428 Second, genome assembly quality can influence TE annotation, particularly in repetitive regions. We
429 showed that LTR annotations are broadly consistent between short-read and long-read assemblies for
430 the same species and that the significant correlation between assembly N50 and genome-wide TE
431 content is primarily driven by satellite sequences rather than LTRs. Further, we identified one case in the
432 Guinea baboon (*Papio papio*) in which local misassembly of the amylase locus substantially affected TE
433 estimates. We resolved this issue by performing targeted local reassembly using the original long reads
434 and re-annotating transposable elements on the corrected locus. Despite these efforts, we cannot entirely
435 rule out similar, undetected local assembly errors in other species, though no additional outliers were
436 observed.

437 Third, the causal direction between LTR enrichment and amylase gene copy number gains remains
438 ambiguous. In SD-rich regions, NAHR can both expand copy number and create a genomic context

439 permissive to additional TE accumulation, making it difficult to disentangle cause from consequence. Our
440 synteny-based evidence indicates that LTR insertions were present in the ancestral segment before
441 duplication, consistent with a seeding role; however, we acknowledge that TE insertions and segmental
442 duplications likely influence each other bidirectionally at this locus.

443 Finally, our regulatory analysis is based on *in silico* TFBS prediction and does not directly measure
444 chromatin accessibility or transcription factor occupancy. The finding that FOXC1 motifs are present in
445 most *AMY* promoters regardless of tissue-specific expression underscores the insufficiency of
446 motif-based prediction alone and highlights the need for functional assays (e.g., ATAC-seq, ChIP-seq,
447 Fiber-seq) to resolve the regulatory architecture of the amylase locus.

448 RESOURCE AVAILABILITY

449 **Lead contact**

450 Requests for further information and resources should be directed to and will be fulfilled by the lead
451 contact, Omer Gokcumen (omergokc@buffalo.edu).

452 **Materials availability**

453 All unique materials generated in this study are available from the lead contact with appropriate
454 institutional approvals.

455 **Data and code availability**

- 456 • Bulk RNA-seq expression data (TPM, GTEx v8) are available from the GTEx open-access portal
457 (https://gtexportal.org/home/downloads/adult-gtex/bulk_tissue_expression). RNA-seq datasets for
458 olive baboon (*Papio anubis*) and rhesus macaque (*Macaca mulatta*) have been deposited at
459 GEO under accession numbers GSE305241 and GSE305255, respectively. All other data
460 reported in this paper are available in the main text and supplementary information.
461 • All original analysis scripts and input files for downstream analyses and visualization have been
462 deposited at Zenodo and are publicly available at <https://doi.org/10.5281/zenodo.16809248> and
463 <https://doi.org/10.5281/zenodo.18689074>. This paper does not report any additional original
464 code.
465 • Any additional information required to reanalyze the data reported in this paper is available from
466 the lead contact upon request

467 ACKNOWLEDGMENTS

468 We thank Leo Speidel, and Luane Landau for carefully reading this manuscript.

469 AUTHOR CONTRIBUTIONS

470 Conceptualization, C.K. and O.G.; methodology, C.K.; Investigation, C.K.; C.K. conducted the
471 bioinformatics analyses; P.P. conducted the ddPCR assay for within-species variation analysis;
472 writing—original draft, C.K. and O.G.; writing—review & editing, C.K., S.R., O.G.; funding acquisition,
473 S.R. and O.G.; resources, O.G.; supervision, O.G.

474 DECLARATION OF INTERESTS

475 The authors declare no competing interests.

476 DECLARATION OF GENERATIVE AI AND AI-ASSISTED TECHNOLOGIES

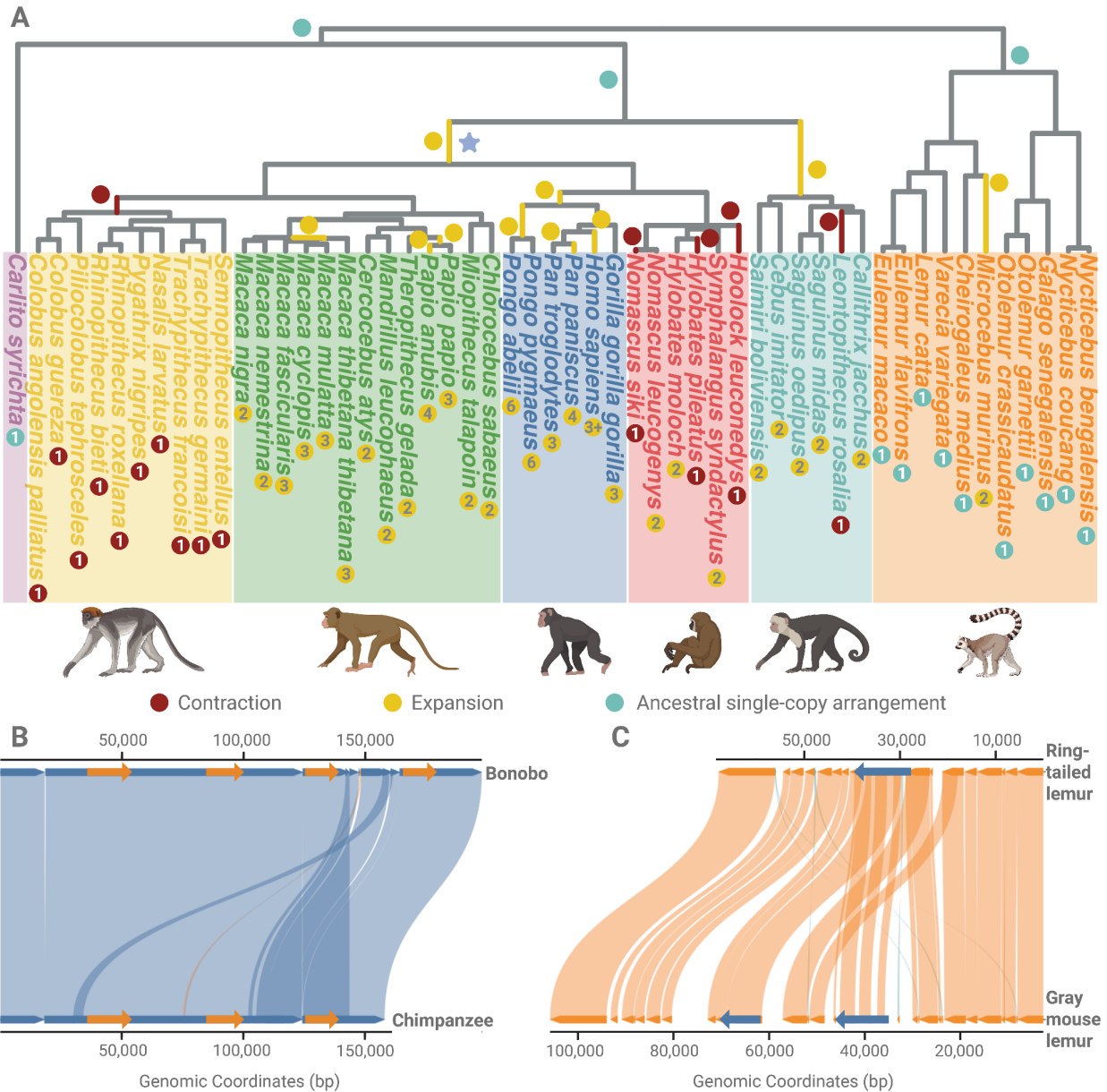
477 During the preparation of this manuscript, the authors used ChatGPT for language and grammar checks.
 478 After using this tool or service, the authors reviewed and edited the content as needed and take full
 479 responsibility for the content of the publication.

480 **SUPPLEMENTAL INFORMATION**

481 **Document S1. Glossary & Figures S1–S14 (pdf file)**

482 **Document S2. Tables S1–S18 (Excel file)**

483 **FIGURE TITLES AND LEGENDS**



484

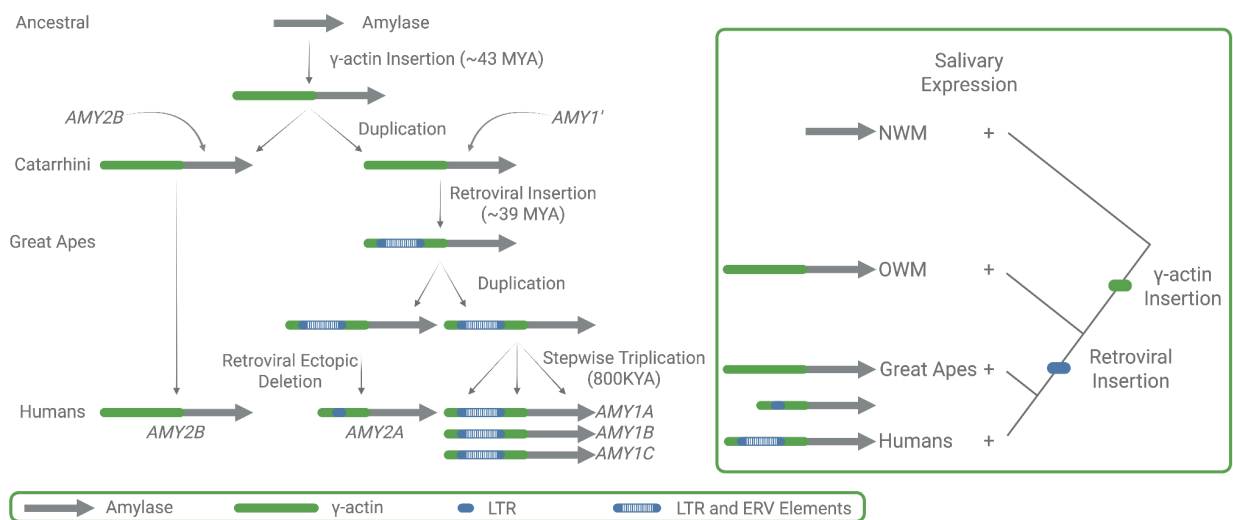
485 **Figure 1. Structural evolutionary history of the amylase locus across primates.**

486 (A) Contractions and expansions in the amylase locus reconstructed from 69 high-quality genomes
 487 representing 53 primate species (Table S2 & S4). Lineages are color-coded by clade: purple for tarsiers
 488 (outgroup), yellow for leaf-eating monkeys, green for Old World monkeys, dark blue for great apes, red
 489 for lesser apes (gibbons), cyan for New World monkeys, and orange for lemurs. Red dots indicate

490 independent contractions in the amylase locus. Yellow dots indicate expansions. Cyan dots mark lineages
 491 retaining the ancestral single-copy configuration, inferred as the ancestral primate state. Independent
 492 contraction events are observed in the gibbon, leaf-eating monkey, and New World monkey genera.
 493 Numbers inside the dots give the total number of *AMY* gene copies detected in each species (per haploid
 494 genome).

495 (B) Synteny comparison between bonobo (*Pan paniscus*) and chimpanzee (*Pan troglodytes*) at the
 496 amylase locus. The copy number increase in bonobo, previously reported¹⁸, is shown here to involve a
 497 chimeric duplication. The 5' flanking region of the duplicated segment resembles the downstream region
 498 of *AMY1*. The internal genic region corresponds to a full duplication of the *AMY1* coding sequence. The 3'
 499 flanking region aligns with the upstream flanking region of *AMY2A*. This chimeric architecture is
 500 consistent with a nonallelic homologous recombination mechanism.

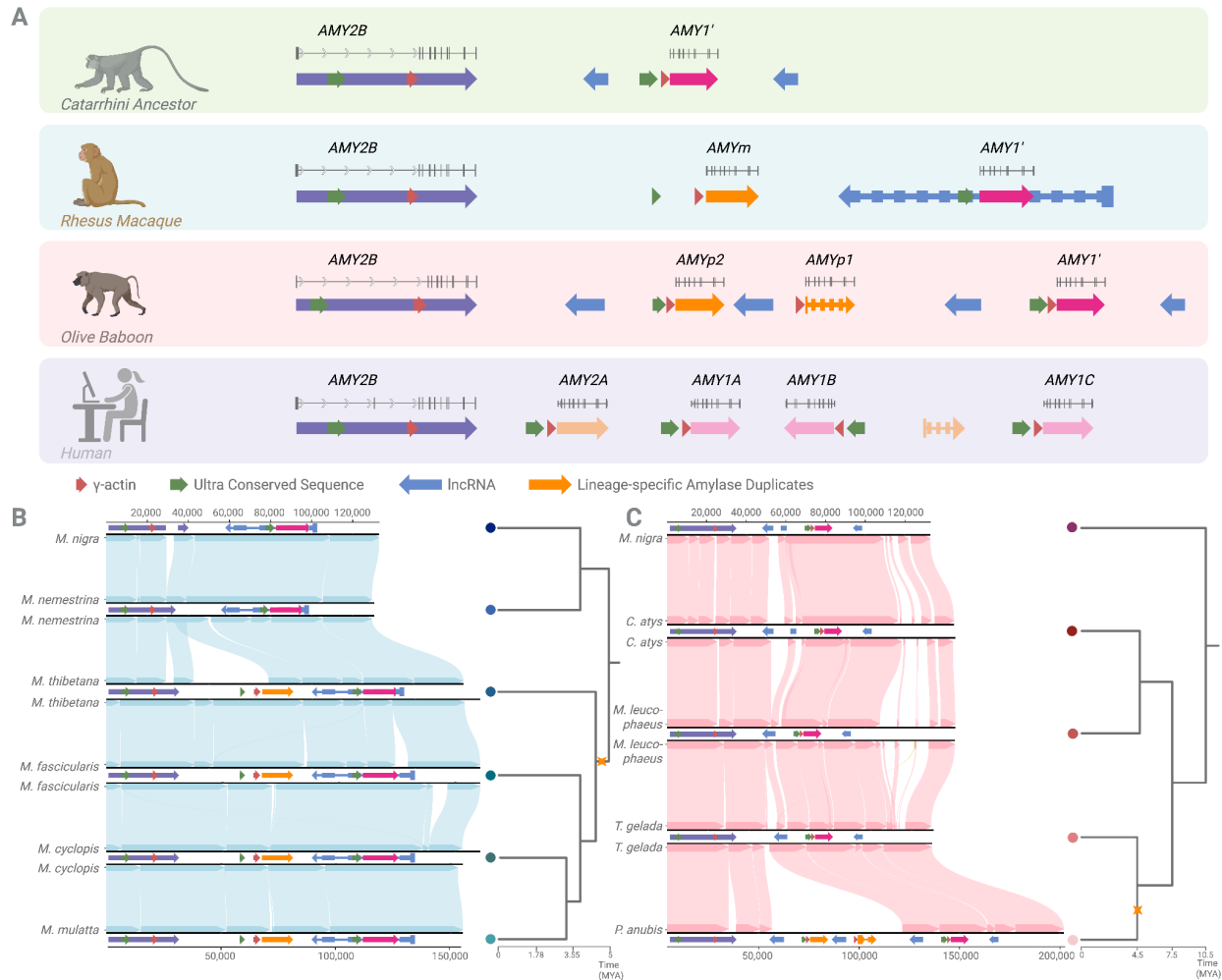
501 (C) Comparison of the amylase locus in ring-tailed lemur (*Lemur catta*) and gray mouse lemur
 502 (*Microcebus murinus*). The gray mouse lemur harbors a tandem duplication of the amylase gene,
 503 including both upstream and downstream flanking regions. This is the only duplication identified among
 504 the lemur species analyzed.



505

506 **Figure 2. Evolutionary origins of the amylase locus in primates.**

507 Model of amylase locus evolution across primates adjusted from Samuelson et al. 1990²⁷. The ancestral
 508 primate genome contained a single-copy amylase gene (*AMY2B*). In the Catarrhini lineage, a γ -actin
 509 insertion occurred upstream of *AMY2B* (~43 MYA), followed by duplication of the γ -actin-*AMY2B*
 510 segment, generating a second gene (*AMY1'*). In great apes, a retroviral (ERV) insertion occurred (~39
 511 MYA), followed by retroviral ectopic deletion in the human lineage, giving rise to *AMY2A*. In humans, a
 512 stepwise triplication of the salivary-expressed amylase gene (*AMY1*) generated *AMY1A*, *AMY1B* and
 513 *AMY1C*. The right panel summarizes the inferred regulatory and structural events, showing independent
 514 gains of salivary expression in New World monkeys, Old World monkeys and great apes, with
 515 subfunctionalization of *AMY1* and *AMY2A* in humans.



516

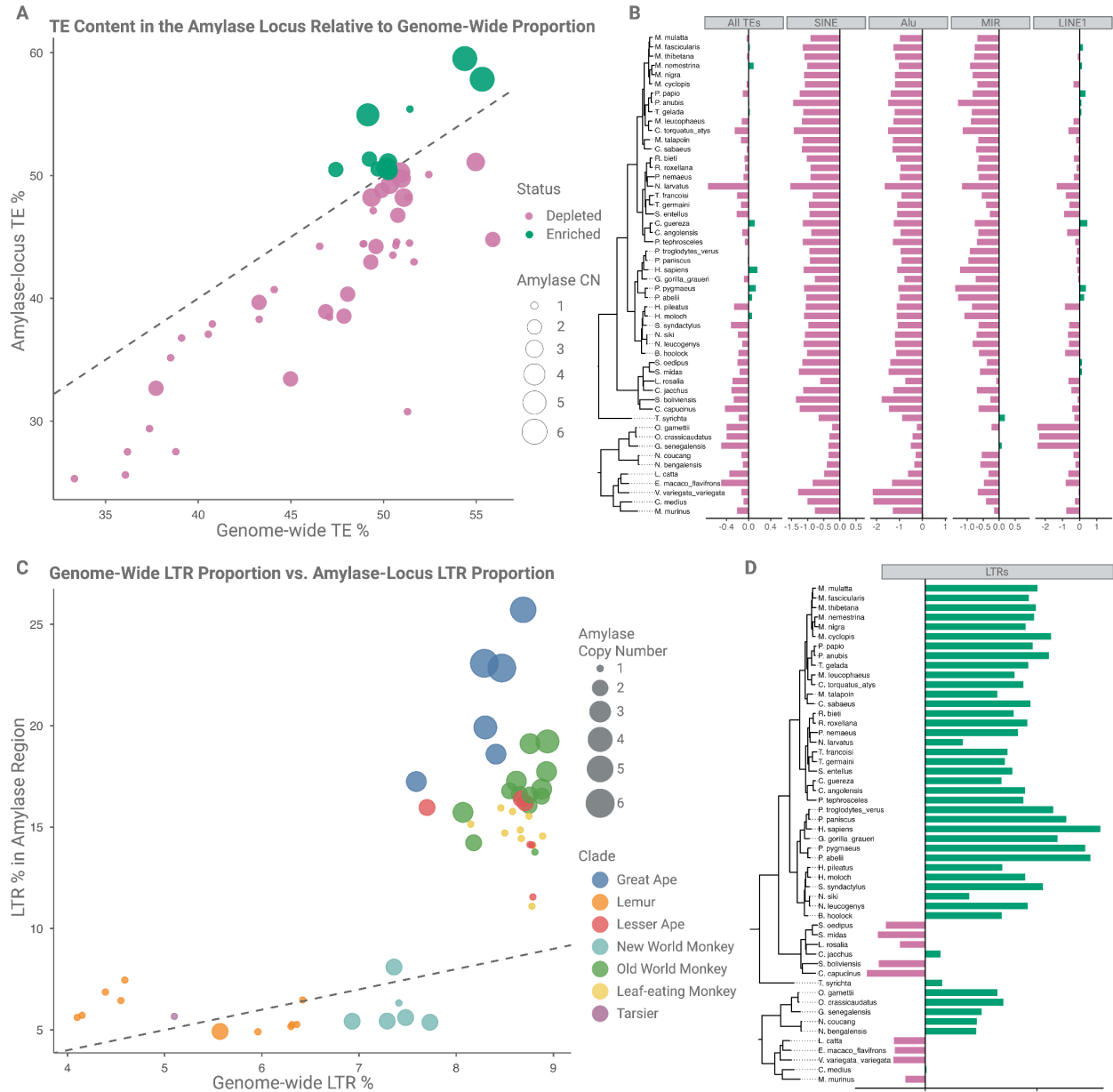
517 **Figure 3. Independent duplication events and structural evolution of the amylase locus in Old**
 518 **World monkeys.**

519 (A) Schematic reconstruction of the amylase locus in the Catarrhini ancestor, rhesus macaque (*Macaca*
 520 *mulatta*), olive baboon (*Papio anubis*), and human reference genome (which is also shown to be the
 521 ancestral human arrangement⁴ (hg38). The Catarrhini ancestor contains two genes: *AMY2B* and *AMY1'*,
 522 the latter derived from a duplication of *AMY2B*. In rhesus macaques (*Macaca mulatta*), the locus includes
 523 *AMY2B*, *AMYm*, and *AMY1'*. In olive baboons (*Papio anubis*), the locus consists of *AMY2B*, *AMYp2* and
 524 *AMYp1*, and *AMY1'*. *AMYp1* is annotated as a pseudogene in NCBI, and is indicated with a dashed
 525 outline. *AMYp1* and *AMYp2* arose from independent duplication events with distinct breakpoints from
 526 those of *AMYm*. In the human reference genome, the locus contains *AMY2B*, *AMY2A* and three *AMY1*
 527 paralogs, *AMY1A*, *AMY1B*, and *AMY1C*. The *AMY2B* genes are one-to-one orthologs across species.
 528 *AMY1'* in macaques, baboons and the Catarrhini ancestor is orthologous among Old World monkeys and
 529 represents the ancestral precursor to human *AMY2A* and *AMY1* genes.

530 (B) Syntenic analysis and phylogenetic context of the rhesus macaque-specific duplication. *AMYm*
 531 (orange) is located between *AMY2B* (purple) and *AMY1'* (pink), and is present only in the fascicularis and
 532 sinica macaque groups, but absent in the silenus group. This distribution dates the duplication to
 533 approximately 4.5-5 million years ago. The gene structure of *AMYm* is identical for coding sequence to
 534 *AMY2B*, but differs on the 5' untranslated region. Breakpoint analysis indicates a nonallelic homologous
 535 recombination (NAHR) mechanism.

536 (C) Syntenic comparison across Papionini species reveals two independent NAHR-mediated duplication
 537 events in olive baboons (*Papio anubis*). The first duplication generated *AMYp1* (orange), and the second
 538 produced *AMYp2* (orange); both events must have occurred after the complete split from Guinea baboon

539 (*Papio papio*), within the ~1.85-million-year window inferred for Papionini divergence. Because of local
 540 misassembly in the amylase locus in Guinea baboon (*Papio papio*), we carried out the synteny analysis
 541 using closely related *Theropithecus gelada* and *Mandrillus leucophaeus* which are members of the
 542 Papionini group, supporting this reconstruction. Species are indicated by Latin names; the corresponding
 543 common names and clade memberships are listed in the **Table S5**.



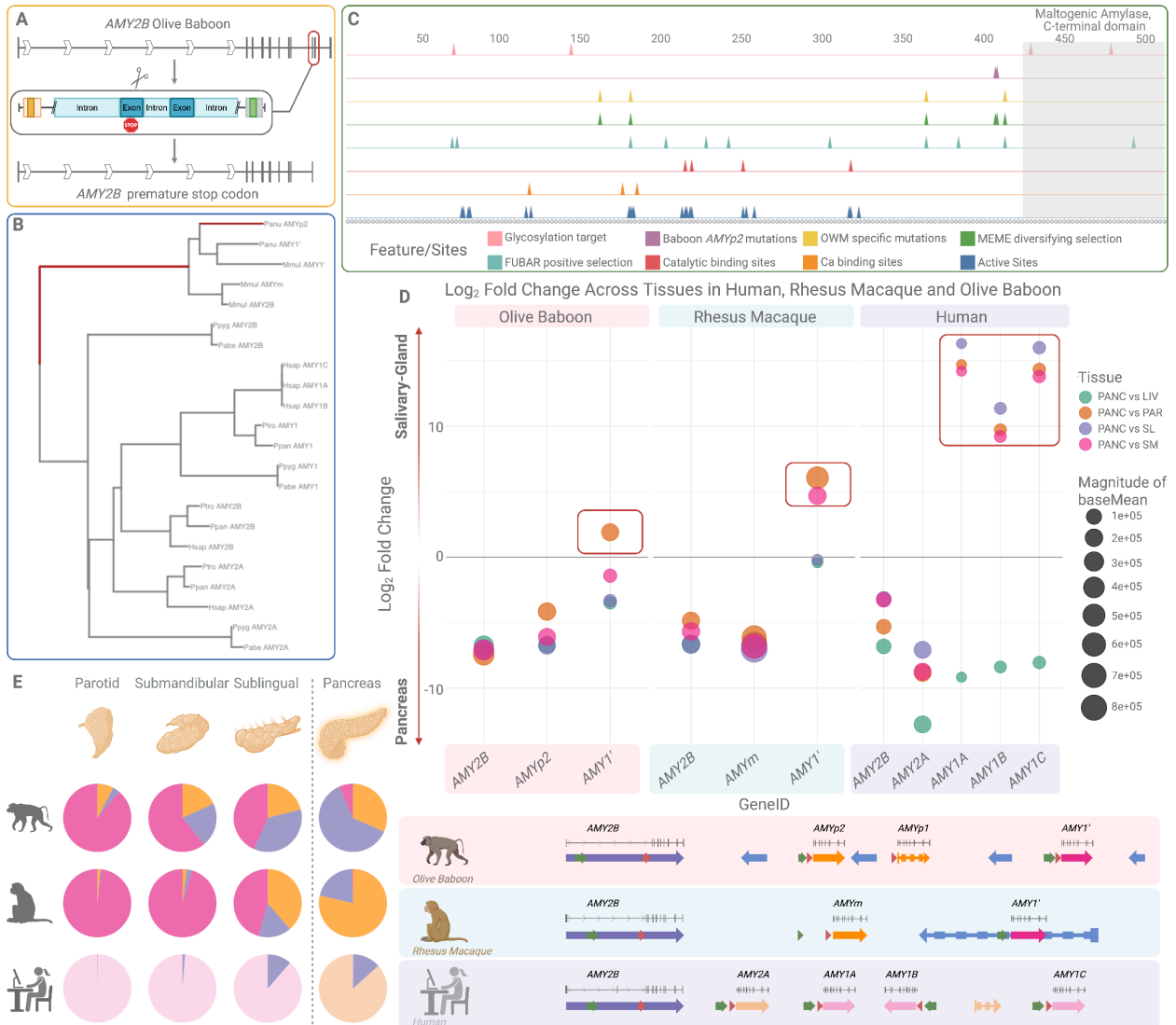
544

545 **Figure 4. The transposable element landscape in the primate amylase locus.**

546 (A) Transposable element (TE) content (given as the total TE in bp/the total length of the locus in bp, see
 547 Methods for more) in the amylase locus relative to genome-wide TE proportion across 53 primate
 548 species. Each point represents a species, with the size of the circle scaled by the respective amylase
 549 copy number and the color indicating enrichment status. The dashed line marks the 1:1 expectation. The
 550 majority of species show depletion of TEs in the amylase locus compared to genome-wide levels.

551 (B) TE family-specific enrichment (\log_2 transformed) in the amylase locus across primates for SINEs,
 552 Alus, MIRs and LINE1s, alongside total TE content (“All TEs”) (Table S6). Bars represent enrichment
 553 (green) or depletion (pink) relative to genome-wide TE representation. Short retrotransposons (e.g. Alus

554 & SINEs) are consistently depleted from the amylase locus, suggesting that this region might have limited
 555 retention of active mobile elements across lineages.
 556 (C) Relationship between genome-wide and amylase locus-specific LTR content across primate species.
 557 Each point represents a species, with the circle size indicating AMY gene copy number and the color
 558 representing phylogenetic clade. The dashed line indicates a 1:1 ratio. Most species show an enrichment
 559 of LTRs at the amylase locus relative to genome-wide levels. (D) Log₂-transformed enrichment and
 560 depletion of LTR content at the amylase locus across 53 primate species, organized by phylogeny. Bars
 561 indicate whether LTR representation in the locus is higher (green) or lower (pink) than genome-wide LTR
 562 proportions. Several species with reported amylase gene duplications, such as olive baboon, exhibit
 563 pronounced LTR enrichment.



564
 565 **Figure 5. Functional divergence and tissue-specific expression of amylase paralogs in olive**
 566 **baboons, rhesus macaques and humans.**
 567 (A) Identification of a premature stop codon in *AMY2B* of olive baboon. A schematic of the disrupted gene
 568 structure (top) shows the location of the nonsense mutation (red box) within the ninth exon.
 569 (B) A maximum likelihood phylogeny of amylase coding sequences from olive baboon, rhesus macaque
 570 and great apes, highlighted branches under significant episodic diversifying selection (dark red) affecting
 571 the lineage leading to Old World monkey paralogs (aBSREL, $p = 0.038$).
 572 (C) Old world monkey and baboon-specific variants and their overlap with predicted functional sites

598 The color of each slice corresponds to the proportion of total expression of the given gene in each tissue,
599 following the same color scheme used for the gene annotation arrows. All three rhesus macaque
600 (*Macaca mulatta*) paralogs (*AMY2B*, *AMYm*, *AMY1'*) share similar TFBS profiles with olive baboons
601 (*Papio anubis*) *AMY1'* and *AMYp2*. In contrast, *AMY2B* in olive baboon and human exhibit distinct TFBS
602 composition. The salivary gland-biased transcription factor FOXC1 is present in most paralogs, but
603 absent from human and olive baboon *AMY2B*, both of which lack salivary expression.

604 REFERENCES

- 605 1. Stern, D.L. (2013). The genetic causes of convergent evolution. *Nat. Rev. Genet.* 14, 751–764.
606 <https://doi.org/10.1038/nrg3483>.
- 607 2. Ferraretti, G., Rill, A., Abondio, P., Smith, K., Ojeda-Granados, C., De Fanti, S., Alberti, M., Izzi, M.,
608 Sherpa, P.T., Cocco, P., et al. (2025). Convergent evolution of complex adaptive traits modulates
609 angiogenesis in high-altitude Andean and Himalayan human populations. *Commun. Biol.* 8, 377.
610 <https://doi.org/10.1038/s42003-025-07813-6>.
- 611 3. Bolognini, D., Halgren, A., Lou, R.N., Raveane, A., Rocha, J.L., Guarracino, A., Soranzo, N., Chin,
612 C.-S., Garrison, E., and Sudmant, P.H. (2024). Recurrent evolution and selection shape structural
613 diversity at the amylase locus. *Nature* 634, 617–625. <https://doi.org/10.1038/s41586-024-07911-1>.
- 614 4. Yilmaz, F., Karageorgiou, C., Kim, K., Pajic, P., Scheer, K., Human Genome Structural Variation
615 Consortium, Beck, C.R., Torregrossa, A.-M., Lee, C., Gokcumen, O., et al. (2024). Reconstruction of
616 the human amylase locus reveals ancient duplications seeding modern-day variation. *Science* 386,
617 eadn0609. <https://doi.org/10.1126/science.adn0609>.
- 618 5. Karageorgiou, C., Gokcumen, O., and Dennis, M.Y. (2024). Deciphering the role of structural
619 variation in human evolution: a functional perspective. *Curr. Opin. Genet. Dev.* 88, 102240.
620 <https://doi.org/10.1016/j.gde.2024.102240>.
- 621 6. Schloissnig, S., Pani, S., Ebler, J., Hain, C., Tsalpalou, V., Söylev, A., Hüther, P., Ashraf, H.,
622 Prodanov, T., Asparuhova, M., et al. (2025). Structural variation in 1,019 diverse humans based on
623 long-read sequencing. *Nature*, 1–11. <https://doi.org/10.1038/s41586-025-09290-7>.
- 624 7. Ohno, S. (2014). *Evolution by Gene Duplication* (Springer).
- 625 8. Moore, R.C., and Purugganan, M.D. (2003). The early stages of duplicate gene evolution. *Proc. Natl.*
626 *Acad. Sci. U. S. A.* 100, 15682–15687. <https://doi.org/10.1073/pnas.2535513100>.
- 627 9. Lan, X., and Pritchard, J.K. (2016). Coregulation of tandem duplicate genes slows evolution of
628 subfunctionalization in mammals. *Science* 352, 1009–1013. <https://doi.org/10.1126/science.aad8411>.
- 629 10. Teufel, A.I., Johnson, M.M., Laurent, J.M., Kachroo, A.H., Marcotte, E.M., and Wilke, C.O. (2019).
630 The many nuanced evolutionary consequences of duplicated genes. *Mol. Biol. Evol.* 36, 304–314.
631 <https://doi.org/10.1093/molbev/msy210>.
- 632 11. Soto, D.C., Uribe-Salazar, J.M., Kaya, G., Valdarrago, R., Sekar, A., Haghani, N.K., Hino, K., La, G.,
633 Mariano, N.A.F., Ingamells, C., et al. (2025). Human-specific gene expansions contribute to brain
634 evolution. *Cell*. <https://doi.org/10.1016/j.cell.2025.06.037>.
- 635 12. Sudmant, P.H., Rausch, T., Gardner, E.J., Handsaker, R.E., Abyzov, A., Huddleston, J., Zhang, Y.,
636 Ye, K., Jun, G., Fritz, M.H.-Y., et al. (2015). An integrated map of structural variation in 2,504 human
637 genomes. *Nature* 526, 75–81. <https://doi.org/10.1038/nature15394>.
- 638 13. Dennis, M.Y., and Eichler, E.E. (2016). Human adaptation and evolution by segmental duplication.
639 *Curr. Opin. Genet. Dev.* 41, 44–52. <https://doi.org/10.1016/j.gde.2016.08.001>.
- 640 14. Roberts, P.J., and Whelan, W.J. (1960). The mechanism of carbohydrase action. 5. Action of human
641 salivary alpha-amylase on amylopectin and glycogen. *Biochem. J.* 76, 246–253.
642 <https://doi.org/10.1042/bj0760246>.
- 643 15. Boehlke, C., Zierau, O., and Hannig, C. (2015). Salivary amylase - The enzyme of unspecialized
644 euryphagous animals. *Arch. Oral Biol.* 60, 1162–1176.
645 <https://doi.org/10.1016/j.archoralbio.2015.05.008>.
- 646 16. Janiak, M.C. (2016). Digestive enzymes of human and nonhuman primates: Digestive enzymes of
647 human and nonhuman primates. *Evol. Anthropol.* 25, 253–266. <https://doi.org/10.1002/evan.21498>.

- 648 17. Pajic, P., Pavlidis, P., Dean, K., Neznanova, L., Romano, R., Garneau, D., Daugherty, E.K., Globig,
649 A., Ruhl, S., and Gokcumen, O. (2019). Independent amylase gene copy number bursts correlate
650 with dietary preferences in mammals. *Elife* 8, e44628. <https://doi.org/10.7554/eLife.44628>.
- 651 18. Perry, G.H., Dominy, N.J., Claw, K.G., Lee, A.S., Fiegler, H., Redon, R., Werner, J., Villanea, F.A.,
652 Mountain, J.L., Misra, R., et al. (2007). Diet and the evolution of human amylase gene copy number
653 variation. *Nat. Genet.* 39, 1256–1260. <https://doi.org/10.1038/ng2123>.
- 654 19. Scheer, K., Landau, L.J.B., Jorgensen, K., Karageorgiou, C., Siao, L., Alkan, C., Morales-Rivera,
655 A.M., Osbourne, C., Garcia, O., Pearson, L., et al. (2025). Adaptive increase of amylase gene copy
656 number in Peruvians driven by potato-rich diets. *bioRxiv*, 2025.03.25.644684.
657 <https://doi.org/10.1101/2025.03.25.644684>.
- 658 20. Soler i Nunez, A., Joly, C., Humbert, C., Chowdhury, A., Green, S.T., Harena, P., Bakrobena, L.,
659 Fomine, F.L.M., Ebbesen, P., Tolesa, Z.G., et al. (2026). Rethinking human AMY1 copy number
660 evolution in light of demographic history. *bioRxiv*, 2026.02.18.706577.
661 <https://doi.org/10.64898/2026.02.18.706577>.
- 662 21. Schibler, U., Pittet, A.C., Young, R.A., Hagenbüchle, O., Tosi, M., Gellman, S., and Wellauer, P.K.
663 (1982). The mouse alpha-amylase multigene family. Sequence organization of members expressed
664 in the pancreas, salivary gland and liver. *J. Mol. Biol.* 155, 247–266.
665 [https://doi.org/10.1016/0022-2836\(82\)90004-3](https://doi.org/10.1016/0022-2836(82)90004-3).
- 666 22. Axelsson, E., Ratnakumar, A., Arendt, M., Maqbool, K., Webster, M., Perloski, M., Liberg, O.,
667 Arnemo, J., Hedhammar, Å., and Lindblad-Toh, K. (2013). The genomic signature of dog
668 domestication reveals adaptation to a starch-rich diet. *Nature* 495, 360–364.
669 <https://doi.org/10.1038/nature11837>.
- 670 23. Behringer, V., Borchers, C., Deschner, T., Möstl, E., Selzer, D., and Hohmann, G. (2013).
671 Measurements of salivary alpha amylase and salivary cortisol in hominoid primates reveal
672 within-species consistency and between-species differences. *PLoS One* 8, e60773.
673 <https://doi.org/10.1371/journal.pone.0060773>.
- 674 24. McGeachin, R.L., and Akin, J.R. (1982). Amylase levels in the tissues and body fluids of several
675 primate species. *Comp. Biochem. Physiol. A Comp. Physiol.* 72, 267–269.
676 [https://doi.org/10.1016/0300-9629\(82\)90045-7](https://doi.org/10.1016/0300-9629(82)90045-7).
- 677 25. Thamadolok, S., Choi, K.-S., Ruhl, L., Schulte, F., Kazim, A.L., Hardt, M., Gokcumen, O., and Ruhl,
678 S. (2020). Human and nonhuman primate lineage-specific footprints in the salivary proteome. *Mol.*
679 *Biol. Evol.* 37, 395–405. <https://doi.org/10.1093/molbev/msz223>.
- 680 26. Samuelson, L.C., Phillips, R.S., and Swanberg, L.J. (1996). Amylase gene structures in primates:
681 retroposon insertions and promoter evolution. *Mol. Biol. Evol.* 13, 767–779.
682 <https://doi.org/10.1093/OXFORDJOURNALS.MOLBEV.A025637>.
- 683 27. Samuelson, L.C., Wiebauer, K., Snow, C.M., and Meisler, M.H. (1990). Retroviral and Pseudogene
684 Insertion Sites Reveal the Lineage of Human Salivary and Pancreatic Amylase Genes from a Single
685 Gene during Primate Evolution. *Molecular and Cellular Biology* 10, 2513–2520.
686 <https://doi.org/10.1128/mcb.10.6.2513-2520.1990>.
- 687 28. Ting, C., Rosenberg, M., Snow, C., Samuelson, L., and Meisler, M. (1992). Endogenous retroviral
688 sequences are required for tissue-specific expression of a human salivary amylase gene. *Genes*
689 *Dev.* 6, 1457–1465. <https://doi.org/10.1101/GAD.6.8.1457>.
- 690 29. Meisler, M.H., and Ting, C.-N. (1993). The Remarkable Evolutionary History of the Human Amylase
691 Genes. *Critical Reviews in Oral Biology & Medicine* 4, 503–509.
692 <https://doi.org/10.1177/10454411930040033501>.
- 693 30. Mandel, A.L., Peyrot des Gachons, C., Plank, K.L., Alarcon, S., and Breslin, P.A.S. (2010). Individual
694 differences in AMY1 gene copy number, salivary α -amylase levels, and the perception of oral starch.
695 *PLoS One* 5, e13352. <https://doi.org/10.1371/journal.pone.0013352>.
- 696 31. Lindsay, S.J., Khajavi, M., Lupski, J.R., and Hurles, M.E. (2006). A chromosomal rearrangement
697 hotspot can be identified from population genetic variation and is coincident with a hotspot for allelic
698 recombination. *Am. J. Hum. Genet.* 79, 890–902. <https://doi.org/10.1086/508709>.
- 699 32. Perry, G.H., Tchinda, J., McGrath, S.D., Zhang, J., Picker, S.R., Cáceres, A.M., lafrate, A.J.,

- 700 Tyler-Smith, C., Scherer, S.W., Eichler, E.E., et al. (2006). Hotspots for copy number variation in
701 chimpanzees and humans. *Proc. Natl. Acad. Sci. U. S. A.* *103*, 8006–8011.
702 <https://doi.org/10.1073/pnas.0602318103>.
- 703 33. Lin, Y.-L., and Gokcumen, O. (2019). Fine-scale characterization of genomic structural variation in
704 the human genome reveals adaptive and biomedically relevant hotspots. *Genome Biol. Evol.* *11*,
705 1136–1151. <https://doi.org/10.1093/gbe/evz058>.
- 706 34. Cooper, E.B., Brent, L., Snyder-Mackler, N., Singh, M., Sengupta, A., Khatiwada, S., Malaivijitnond,
707 S., Hai, Z.Q., and Higham, J. (2022). The natural history of model organisms: the rhesus macaque
708 as a success story of the Anthropocene. *Elife* *11*. <https://doi.org/10.7554/eLife.78169>.
- 709 35. Rogers, J., Raveendran, M., Harris, R.A., Mailund, T., Leppälä, K., Athanasiadis, G., Schierup, M.H.,
710 Cheng, J., Munch, K., Walker, J.A., et al. (2019). The comparative genomics and complex population
711 history of *Papio* baboons. *Sci. Adv.* *5*, eaau6947. <https://doi.org/10.1126/sciadv.aau6947>.
- 712 36. Roos, C., Knauf, S., Chuma, I.S., Maille, A., Callou, C., Sabin, R., Portela Miguez, R., and Zinner, D.
713 (2021). New mitogenomic lineages in *Papio* baboons and their phylogeographic implications. *Am. J.*
714 *Phys. Anthropol.* *174*, 407–417. <https://doi.org/10.1002/ajpa.24186>.
- 715 37. Gu, W., Zhang, F., and Lupski, J.R. (2008). Mechanisms for human genomic rearrangements.
716 *Pathogenetics* *1*, 4. <https://doi.org/10.1186/1755-8417-1-4>.
- 717 38. Lieber, M.R. (2008). The mechanism of human nonhomologous DNA end joining. *J. Biol. Chem.*
718 *283*, 1–5. <https://doi.org/10.1074/jbc.R700039200>.
- 719 39. Zhang, F., Khajavi, M., Connolly, A.M., Towne, C.F., Batish, S.D., and Lupski, J.R. (2009). The DNA
720 replication FoSTeS/MMBIR mechanism can generate genomic, genic and exonic complex
721 rearrangements in humans. *Nat. Genet.* *41*, 849–853. <https://doi.org/10.1038/ng.399>.
- 722 40. Reiter, L.T., Hastings, P.J., Nelis, E., De Jonghe, P., Van Broeckhoven, C., and Lupski, J.R. (1998).
723 Human meiotic recombination products revealed by sequencing a hotspot for homologous strand
724 exchange in multiple HNPP deletion patients. *Am. J. Hum. Genet.* *62*, 1023–1033.
725 <https://doi.org/10.1086/301827>.
- 726 41. Aquadro, C.F., Weaver, A.L., Schaeffer, S.W., and Anderson, W.W. (1991). Molecular evolution of
727 inversions in *Drosophila pseudoobscura*: the amylase gene region. *Proc. Natl. Acad. Sci. U. S. A.*
728 *88*, 305–309. <https://doi.org/10.1073/pnas.88.1.305>.
- 729 42. Staubach, F., Lorenc, A., Messer, P.W., Tang, K., Petrov, D.A., and Tautz, D. (2012). Genome
730 patterns of selection and introgression of haplotypes in natural populations of the house mouse (*Mus*
731 *musculus*). *PLoS Genet.* *8*, e1002891. <https://doi.org/10.1371/journal.pgen.1002891>.
- 732 43. Johnson, M.E., National Institute of Health Intramural Sequencing Center Comparative Sequencing
733 Program, Cheng, Z., Morrison, V.A., Scherer, S., Ventura, M., Gibbs, R.A., Green, E.D., and Eichler,
734 E.E. (2006). Recurrent duplication-driven transposition of DNA during hominoid evolution. *Proc. Natl.*
735 *Acad. Sci. U. S. A.* *103*, 17626–17631. <https://doi.org/10.1073/pnas.0605426103>.
- 736 44. Balachandran, P., Walawalkar, I.A., Flores, J.I., Dayton, J.N., Audano, P.A., and Beck, C.R. (2022).
737 Transposable element-mediated rearrangements are prevalent in human genomes. *Nat. Commun.*
738 *13*, 7115. <https://doi.org/10.1038/s41467-022-34810-8>.
- 739 45. Kamitaki, N., Handsaker, R.E., Hujoel, M.L.A., Mukamel, R.E., Usher, C.L., McCarroll, S.A., and Loh,
740 P.-R. (2026). Human and bacterial genetic variation shape oral microbiomes and health. *Nature*,
741 1–11. <https://doi.org/10.1038/s41586-025-10037-7>.
- 742 46. Takuno, S., Nishio, T., Satta, Y., and Innan, H. (2008). Preservation of a pseudogene by gene
743 conversion and diversifying selection. *Genetics* *180*, 517–531.
744 <https://doi.org/10.1534/genetics.108.091918>.
- 745 47. Velová, H., Gutowska-Ding, M.W., Burt, D.W., and Vinkler, M. (2018). Toll-like receptor evolution in
746 birds: Gene duplication, pseudogenization, and diversifying selection. *Mol. Biol. Evol.* *35*,
747 2170–2184. <https://doi.org/10.1093/molbev/msy119>.
- 748 48. Wang, J., Chitsaz, F., Derbyshire, M.K., Gonzales, N.R., Gwadz, M., Lu, S., Marchler, G.H., Song,
749 J.S., Thanki, N., Yamashita, R.A., et al. (2023). The conserved domain database in 2023. *Nucleic*
750 *Acids Res.* *51*, D384–D388. <https://doi.org/10.1093/nar/gkac1096>.

- 751 49. Saitou, M., Gaylord, E.A., Xu, E., May, A.J., Neznanova, L., Nathan, S., Grawe, A., Chang, J., Ryan,
752 W., Ruhl, S., et al. (2020). Functional specialization of human salivary glands and origins of proteins
753 intrinsic to human saliva. *Cell Rep.* 33, 108402. <https://doi.org/10.1016/j.celrep.2020.108402>.
- 754 50. Bray, N.L., Pimentel, H., Melsted, P., and Pachter, L. (2016). Near-optimal probabilistic RNA-seq
755 quantification. *Nat. Biotechnol.* 34, 525–527. <https://doi.org/10.1038/nbt.3519>.
- 756 51. Michael, D.G., Pranzatelli, T.J.F., Warner, B.M., Yin, H., and Chiorini, J.A. (2019). Integrated
757 epigenetic mapping of human and mouse salivary gene regulation. *J. Dent. Res.* 98, 209–217.
758 <https://doi.org/10.1177/0022034518806518>.
- 759 52. Kamitaki, N., Handsaker, R.E., Hujuel, M.L.A., Mukamel, R.E., Usher, C.L., McCarroll, S.A., and Loh,
760 P.-R. (2025). Human and bacterial genetic variation shape oral microbiomes and health. medRxiv,
761 2025.03.31.25324952. <https://doi.org/10.1101/2025.03.31.25324952>.
- 762 53. Ponting, C.P. (2017). Biological function in the twilight zone of sequence conservation. *BMC Biol.* 15,
763 71. <https://doi.org/10.1186/s12915-017-0411-5>.
- 764 54. Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local alignment
765 search tool. *J. Mol. Biol.* 215, 403–410. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2).
- 766 55. Danecek, P., Bonfield, J.K., Liddle, J., Marshall, J., Ohan, V., Pollard, M.O., Whitwham, A., Keane,
767 T., McCarthy, S.A., Davies, R.M., et al. (2021). Twelve years of SAMtools and BCFtools.
768 *Gigascience* 10, giab008. <https://doi.org/10.1093/gigascience/giab008>.
- 769 56. Kurtz, S., Phillippy, A., Delcher, A.L., Smoot, M., Shumway, M., Antonescu, C., and Salzberg, S.L.
770 (2004). Versatile and open software for comparing large genomes. *Genome Biol.* 5, R12.
771 <https://doi.org/10.1186/gb-2004-5-2-r12>.
- 772 57. Kielbasa, S.M., Wan, R., Sato, K., Horton, P., and Frith, M.C. (2011). Adaptive seeds tame genomic
773 sequence comparison. *Genome Res.* 21, 487–493. <https://doi.org/10.1101/gr.113985.110>.
- 774 58. Porubsky, D., Guitart, X., Yoo, D., Dishuck, P.C., Harvey, W.T., and Eichler, E.E. (2025). SVbyEye: a
775 visual tool to characterize structural variation among whole-genome assemblies. *Bioinformatics* 41,
776 btaf332. <https://doi.org/10.1093/bioinformatics/btaf332>.
- 777 59. Slater, G.S.C., and Birney, E. (2005). Automated generation of heuristics for biological sequence
778 comparison. *BMC Bioinformatics* 6, 31. <https://doi.org/10.1186/1471-2105-6-31>.
- 779 60. Schliep, K.P. (2011). phangorn: phylogenetic analysis in R. *Bioinformatics* 27, 592–593.
780 <https://doi.org/10.1093/bioinformatics/btq706>.
- 781 61. Paradis, E., and Schliep, K. (2019). ape 5.0: an environment for modern phylogenetics and
782 evolutionary analyses in R. *Bioinformatics* 35, 526–528.
783 <https://doi.org/10.1093/bioinformatics/bty633>.
- 784 62. Išerić, H., Alkan, C., Hach, F., and Numanagić, I. (2022). Fast characterization of segmental
785 duplication structure in multiple genome assemblies. *Algorithms Mol. Biol.* 17, 4.
786 <https://doi.org/10.1186/s13015-022-00210-2>.
- 787 63. Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34,
788 3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>.
- 789 64. Kolmogorov, M., Bickhart, D.M., Behsaz, B., Gurevich, A., Rayko, M., Shin, S.B., Kuhn, K., Yuan, J.,
790 Polevikov, E., Smith, T.P.L., et al. (2020). metaFlye: scalable long-read metagenome assembly using
791 repeat graphs. *Nat. Methods* 17, 1103–1110. <https://doi.org/10.1038/s41592-020-00971-x>.
- 792 65. Cheng, H., Asri, M., Lucas, J., Koren, S., and Li, H. (2024). Scalable telomere-to-telomere assembly
793 for diploid and polyploid genomes with double graph. *Nat. Methods* 21, 967–970.
794 <https://doi.org/10.1038/s41592-024-02269-8>.
- 795 66. Wingett, S.W., and Andrews, S. (2018). FastQ Screen: A tool for multi-genome mapping and quality
796 control. *F1000Res.* 7, 1338. <https://doi.org/10.12688/f1000research.15931.2>.
- 797 67. Ewels, P., Magnusson, M., Lundin, S., and Käller, M. (2016). MultiQC: summarize analysis results for
798 multiple tools and samples in a single report. *Bioinformatics* 32, 3047–3048.
799 <https://doi.org/10.1093/bioinformatics/btw354>.
- 800 68. Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads.
801 *EMBnet J.* 17, 10. <https://doi.org/10.14806/ej.17.1.200>.

- 802 69. Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion
803 for RNA-seq data with DESeq2. *Genome Biol.* 15, 550. <https://doi.org/10.1186/s13059-014-0550-8>.
- 804 70. Kim, D., Paggi, J.M., Park, C., Bennett, C., and Salzberg, S.L. (2019). Graph-based genome
805 alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* 37, 907–915.
806 <https://doi.org/10.1038/s41587-019-0201-4>.
- 807 71. Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler
808 transform. *Bioinformatics* 25, 1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>.
- 809 72. Katoh, K., and Standley, D.M. (2013). MAFFT multiple sequence alignment software version 7:
810 improvements in performance and usability. *Mol. Biol. Evol.* 30, 772–780.
811 <https://doi.org/10.1093/molbev/mst010>.
- 812 73. Suyama, M., Torrents, D., and Bork, P. (2006). PAL2NAL: robust conversion of protein sequence
813 alignments into the corresponding codon alignments. *Nucleic Acids Res.* 34, W609–W612.
814 <https://doi.org/10.1093/nar/gkl315>.
- 815 74. Nguyen, L.-T., Schmidt, H.A., von Haeseler, A., and Minh, B.Q. (2015). IQ-TREE: a fast and effective
816 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* 32, 268–274.
817 <https://doi.org/10.1093/molbev/msu300>.
- 818 75. Smith, M.D., Wertheim, J.O., Weaver, S., Murrell, B., Scheffler, K., and Kosakovsky Pond, S.L.
819 (2015). Less is more: an adaptive branch-site random effects model for efficient detection of episodic
820 diversifying selection. *Mol. Biol. Evol.* 32, 1342–1353. <https://doi.org/10.1093/molbev/msv022>.
- 821 76. Murrell, B., Wertheim, J.O., Moola, S., Weighill, T., Scheffler, K., and Kosakovsky Pond, S.L. (2012).
822 Detecting individual sites subject to episodic diversifying selection. *PLoS Genet.* 8, e1002764.
823 <https://doi.org/10.1371/journal.pgen.1002764>.
- 824 77. Murrell, B., Moola, S., Mabona, A., Weighill, T., Sheward, D., Kosakovsky Pond, S.L., and Scheffler,
825 K. (2013). FUBAR: a fast, unconstrained bayesian approximation for inferring selection. *Mol. Biol.*
826 *Evol.* 30, 1196–1205. <https://doi.org/10.1093/molbev/mst030>.
- 827 78. Wertheim, J.O., Murrell, B., Smith, M.D., Kosakovsky Pond, S.L., and Scheffler, K. (2015). RELAX:
828 detecting relaxed selection in a phylogenetic framework. *Mol. Biol. Evol.* 32, 820–832.
829 <https://doi.org/10.1093/molbev/msu400>.
- 830 79. Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K.,
831 Bates, R., Žídek, A., Potapenko, A., et al. (2021). Highly accurate protein structure prediction with
832 AlphaFold. *Nature* 596, 583–589. <https://doi.org/10.1038/s41586-021-03819-2>.
- 833 80. Ramasubbu, N., Paloth, V., Luo, Y., Brayer, G.D., and Levine, M.J. (1996). Structure of human
834 salivary α -amylase at 1.6 Å resolution: Implications for its role in the oral cavity. *Acta Crystallogr. D*
835 *Biol. Crystallogr.* 52, 435–446. <https://doi.org/10.1107/S09074444995014119>.
- 836 81. Ramasubbu, N., Rangunath, C., and Mishra, P.J. (2003). Probing the role of a mobile loop in
837 substrate binding and enzyme activity of human salivary amylase. *J. Mol. Biol.* 325, 1061–1076.
838 [https://doi.org/10.1016/s0022-2836\(02\)01326-8](https://doi.org/10.1016/s0022-2836(02)01326-8).
- 839 82. Gupta, R., and Brunak, S. (2002). Prediction of glycosylation across the human proteome and the
840 correlation to protein function. *Pac. Symp. Biocomput.*, 310–322.
- 841 83. Dreos, R., Ambrosini, G., Périer, R.C., and Bucher, P. (2015). The Eukaryotic Promoter Database:
842 expansion of EPDnew and new promoter analysis tools. *Nucleic Acids Res.* 43, D92–D96.
843 <https://doi.org/10.1093/nar/gku1111>.
- 844 84. Bailey, T.L., and Elkan, C. (1994). Fitting a mixture model by expectation maximization to discover
845 motifs in biopolymers. *Proc. Int. Conf. Intell. Syst. Mol. Biol.* 2, 28–36.
- 846 85. Gupta, S., Stamatoyannopoulos, J.A., Bailey, T.L., and Noble, W.S. (2007). Quantifying similarity
847 between motifs. *Genome Biol.* 8, R24. <https://doi.org/10.1186/gb-2007-8-2-r24>.
- 848 86. Rauluseviciute, I., Riudavets-Puig, R., Blanc-Mathieu, R., Castro-Mondragon, J.A., Ferenc, K.,
849 Kumar, V., Lemma, R.B., Lucas, J., Chèneby, J., Baranasic, D., et al. (2024). JASPAR 2024: 20th
850 anniversary of the open-access database of transcription factor binding profiles. *Nucleic Acids Res.*
851 52, D174–D182. <https://doi.org/10.1093/nar/gkad1059>.
- 852 87. Noguchi, S., Arakawa, T., Fukuda, S., Furuno, M., Hasegawa, A., Hori, F., Ishikawa-Kato, S., Kaida,

853 K., Kaiho, A., Kanamori-Katayama, M., et al. (2017). FANTOM5 CAGE profiles of human and mouse
854 samples. *Sci. Data* 4, 170112. <https://doi.org/10.1038/sdata.2017.112>.
855 88. Grant, C.E., Bailey, T.L., and Noble, W.S. (2011). FIMO: scanning for occurrences of a given motif.
856 *Bioinformatics* 27, 1017–1018. <https://doi.org/10.1093/bioinformatics/btr064>.

857

858 EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS

859 Macaque and Baboon salivary glands samples were obtained from Texas Biomedical Institute through
860 their established IACUC protocols and by expert veterinary staff. The samples were obtained
861 opportunistically during planned euthanasia from animals that are 7 to 18 years old (**Table S8**).

862 METHOD DETAILS

863 *Primate Genome Assemblies and Locus Extraction*

864 To capture the phylogenetic diversity of the amylase locus across primates, we analyzed 244 high-quality
865 genome assemblies representing 222 species (**Table S1**). To assess copy number variation and
866 structural organization, we first delineated the locus using two flanking, non copy-number-variable anchor
867 sequences: one upstream (5') of the *AMY2B* gene corresponding to the *RNPC3* genic region, and one
868 downstream (3') of the *AMY1C* gene. Both sequences were derived from the human GRCh38 reference
869 genome and used as queries in blastn (v. 2.14.1+)⁵⁴ searches against each assembly.

870 We automated this analysis using a custom SLURM-based Bash pipeline. The script identified
871 assemblies where both anchors mapped to a single contig or scaffold, extracted the intervening locus
872 using samtools faidx (v. 1.22)⁵⁵ and generated reverse complements with seqtk (v. 1.5-r133;
873 <https://github.com/lh3/seqtk>) for loci located on the reverse strand. Only assemblies in which the full
874 amylase locus was contained within a single contig were retained for downstream analyses; loci
875 fragmented across multiple scaffolds were excluded but recorded. The queried sequences and the
876 SLURM-based Bash pipeline have been deposited in Zenodo; see Data Availability.

877 *Synteny Analyses and Gene Annotation*

878 To characterize the structural composition and assess synteny of the amylase locus across primates, we
879 analyzed the 70 assemblies in which the complete locus was successfully extracted as a single
880 contiguous sequence. Each locus was aligned to the human H1^a.1 haplotype, which carries a single
881 *AMY1* gene and represents the ancestral configuration in great apes, using NUCmer (v. 3.1)⁵⁶ and
882 LAST⁵⁷. Dotplots were then generated using mummerplot (v. 3.5)⁵⁶ and visually inspected to identify
883 structural rearrangements relative to the human reference.

884 To further investigate patterns of structural evolution both within and between genera, we performed
885 additional pairwise alignments among species of the same genus, as well as representative comparisons
886 across genera. These alignments were again generated with NUCmer (v. 3.1) and visualized using
887 dotplots and miropeats-style plots. For the latter, we used custom scripts to convert NUCmer coordinate
888 files into PAF format and visualized the alignments using the R package SVbyEye (v. 0.99.0)⁵⁸.

889 To annotate amylase gene copies across primates, we employed a multi-step approach combining
890 reference annotations, CDS mapping, and manual curation. For each genome assembly with a
891 contiguous extracted amylase locus, we first retrieved available NCBI gene annotations for that species
892 and used exonerate (v. 2.4.0)⁵⁹ to map the corresponding coding sequences (CDS) to the extracted
893 regions. In assemblies lacking gene annotations, we utilized CDS from closely related sister species and
894 mapped them to the respective genomes using the same exonerate-based pipeline, which is well-suited
895 for aligning sequences in the presence of sequence divergence and detecting partial exons.

896 To identify putative orthologs and resolve paralogous relationships, we implemented a recursive
897 reciprocal BLAST approach. Putative amylase gene models from each species were queried against one

898 another using blastn (v. 2.14.1+) and reciprocal searches were used to establish one-to-one or
899 one-to-many orthology/paralogy relationships. In a one-to-one orthology relationship, reciprocal BLAST
900 searches yield the same gene pair as the top hit in both directions e.g. gene X from species 1 identifies
901 gene Y from species 2, and gene Y reciprocally returns gene X. By contrast, a one-to-many pattern arises
902 when duplication has occurred in only one lineage, and this can appear in two complementary ways. If
903 species 2 experienced the duplication, a single gene in species 1 matches several genes in species 2,
904 and each of those paralogs reciprocally lists the original query gene as its best hit. Conversely, if the
905 duplication occurred in species 1, multiple paralogs in that species all identify the same ortholog in
906 species 2 as their top hit, while the lone gene in species 2 reciprocally aligns most strongly to just one,
907 commonly the most conserved, of the copies in species 1. Either configuration reveals a one-to-many
908 orthology generated by a lineage-specific duplication event. This approach allowed us to confidently
909 distinguish ancestral copies (e.g. *AMY2B*) from lineage-specific duplications.

910 Lastly, in olive baboons (*Papio anubis*), we detected a premature stop codon within the *AMY2B* gene,
911 suggesting likely pseudogenization. This mutation is absent from the NCBI annotation and was not
912 present in the CDS of Guinea baboon (*Papio papio*), but was consistently identified in both haploid
913 assemblies of olive baboons.

914 **Transposable Element Annotation and Quantification**

915 All 53 primate assemblies that yielded a contiguous amylase locus were analyzed using RepeatMasker
916 (v. 4.1.5; <http://repeatmasker.org>) with the -species primates flag and default parameters. This approach
917 applies a uniform primate-specific repeat library to every genome, avoiding the inconsistencies that would
918 arise from mixing lineage-specific libraries. However, this choice can systematically underestimate
919 species-restricted transposable element (TE) families absent from the reference library, an acknowledged
920 limitation of our approach. RepeatMasker output files (.tbl and .out) were parsed to calculate
921 genome-wide TE content for each assembly. Assembly quality introduces additional bias to these
922 estimates, as chromosome-level assemblies often recover TE-rich centromeric and telomeric regions that
923 incomplete or nearly-complete assemblies commonly miss, thereby inflating total TE proportions.

924 To analyze TE content within the amylase locus, we applied RepeatMasker separately to the extracted
925 amylase sequences from each of the 53 species. These sequences were bounded by non
926 copy-number-variant *RNPC3* 5' and *AMY1C* 3' flanks (detailed in the previous section) and restricted to
927 loci assembled as single, gap-free contigs/scaffolds, as a prerequisite to avoid potential segmental
928 duplication collapse. Masking isolated locus sequences, rather than extracting annotations from
929 whole-genome GFF3 files, prevented coordinate errors and ensured that every species was compared
930 across identical boundaries.

931 TE abundance was quantified as the proportion of masked bases (bp masked / region length) for both
932 genome-wide and locus-specific analyses. This normalization accounts for length variation introduced by
933 gene duplications. All downstream statistical analyses and visualizations were conducted in R (v. 4.3.3).
934 For each species, we compared the percentage of masked sequence in the amylase locus to its
935 genome-wide percentage. For each TE family, we calculated log₂-enrichment, defined as log₂(% in
936 locus)/(% genome-wide). This metric was then merged with amylase copy number estimates and clade
937 assignments. Enrichment or depletion was assessed with a two-sided Wilcoxon signed-rank test across
938 the 53 species in which the locus was contiguous; family-specific analyses (e.g. Alu, LTR, DNA
939 transposons) were performed analogously. P-values were adjusted for multiple testing using the
940 Benjamini-Hochberg procedure, considering FDR<0.05 significant.

941 We assessed the potential confounding effect of assembly contiguity (scaffold N50/Assembly N50) on TE
942 content estimates. Assembly N50 showed a strong positive correlation with genome-wide TE content
943 (Spearman $\rho=0.554$, $P<0.001$; $n=53$). However, this relationship was primarily driven by assembly length.
944 After controlling for total assembly length, the association between N50 and TE content remained
945 significant but was substantially reduced (partial correlation $\rho=0.327$, $P=0.018$), indicating that assembly
946 quality has a modest independent effect on TE detection beyond simple genome size effects.

947 To directly evaluate the impact of sequencing technology on TE estimation, we leveraged two species for
948 which both short-read-based and long-read-based assemblies are available (*Semnopithecus entellus* and
949 *Macaca nemestrina*). For each species, we ran RepeatMasker on both assemblies with identical
950 parameters and quantified the genome-wide proportion of sequence annotated as major TE classes
951 (LINEs, SINEs, LTR retrotransposons, DNA transposons, satellites/low-complexity). Across both species,
952 total TE content and the relative contributions of LINEs, SINEs and LTR retrotransposons were highly
953 similar between short-read and long-read assemblies, whereas the long-read assembly in *Macaca*
954 *nemestrina* recovered a higher fraction of sequence annotated as satellites and related tandem repeats,
955 consistent with improved representation of centromeric and pericentromeric regions in more contiguous
956 genomes (**Table S10 & S11**).

957 To test whether TE abundance predicts amylase gene copy number while accounting for shared ancestry,
958 we fit phylogenetic generalized least-squares (PGLS) models in R (nlme 3.1 and caper 1.0.1). A
959 ultrametric primate tree pruned to the 53 focal species was obtained from TimeTree (accessed January
960 2025). The response variable was amylase gene copy number while the predictor was locus-specific TE
961 proportion (total or by family). Pagel's λ was estimated by maximum likelihood and retained if significantly
962 different from zero. Model fit was evaluated by R^2 and residual diagnostics. All scripts for TE analysis,
963 including those used to reproduce Figure 4, are available in a Zenodo repository (see Data Availability)

964 **LTR Orthogroup Analysis and Ancestral-State Reconstruction**

965 To assess the evolutionary dynamics of individual LTR insertions at the amylase locus across primates,
966 we identified orthologous insertions using a reciprocal flanking-sequence approach. For each LTR
967 element annotated by RepeatMasker within the extracted amylase locus of each species, we extracted
968 1kb flanking sequences on both sides and used reciprocal BLAST to identify orthologous insertions
969 across species. Elements sharing orthologous flanking context were clustered into orthogroups. A binary
970 presence/absence matrix (214 orthogroups x 53 species) was then constructed and analyzed on the
971 pruned primate phylogeny.

972 We computed parsimony steps per orthogroup using phangorn v.2.12⁶⁰ and retained orthogroups
973 requiring ≤ 10 changes across the tree which included all trees ($n=214$). For each retained orthogroup, we
974 reconstructed ancestral states using maximum-likelihood discrete-character analysis (ace function in ape
975 v.5.8)⁶¹ under an equal-rates (ER) model. MAP ancestral states were assigned to internal nodes, and
976 branch-wise gains and losses were tallied by comparing parent-child state pairs across all edges and
977 orthogroups. Results were visualized by painting branch colors proportional to the total number of LTR
978 gains or losses on the primate phylogeny (**Figure S9 & S10**). All scripts are deposited in the Zenodo
979 repository (see Data Availability).

980 **Structural Variant Inference and Duplication Mechanism Analysis**

981 To investigate the mechanisms underlying the duplication in the amylase locus, we extended our
982 annotation of the region beyond the existing annotations available on NCBI. We utilized BISER (v. 1.4)⁶²
983 to identify segmental duplications in the olive baboon amylase locus, analyzing the masked and
984 unmasked genome to capture all duplicated units, irrespective of functional annotation. This approach led
985 to the identification of 43 distinct duplicated segments (**Table S12**). Additionally, we identified four lncRNA
986 sequences within the locus, positioned downstream of each gene with reverse orientations relative to the
987 genes. Using self-alignments, dotplots, and BLAST searches, we partitioned the amylase locus into four
988 distinct segments, each corresponding to a single duplicon (**Figure S13A**). Each segment includes the
989 coding sequence of an amylase gene and extends to the terminal sequence of an adjacent lncRNA. This
990 partition enables a systematic examination of duplicated regions across the locus.

991 We then compared the four distinct segments to one another using pairwise NUCmer alignments and
992 visualized these comparisons with dotplots. Query coverage and sequence similarity metrics allowed us
993 to infer the order of duplication events. A closer examination of the nucleotide sequences further informed
994 us of the duplication mechanisms involved. Namely, the third segment (containing the amylase gene
995 *AMYp1*) appears to be a composite of the fourth segment (encompassing *AMY1'*) and the first segment
996 (encompassing *AMY2B*) (**Figure S13B**). Meanwhile, the second segment is nearly identical to the third,

997 with the only difference being a ~10 kb deletion within the third segment (**Figure S13C**). Because this
998 interval is bounded by long stretches of Ns in the olive baboon assembly, we cannot determine whether it
999 represents a true biological deletion or a local misassembly/scaffolding artifact, and we therefore do not
1000 interpret this feature further. Taken together, the sequence-similarity patterns indicate that two non-allelic
1001 homologous recombination (NAHR) events duplicated the ancestral *AMY2B* and *AMY1'* blocks, yielding
1002 the present-day amylase locus structure in olive baboons (**Figure S13**).

1003

1004 To delineate precisely the novel duplications in the olive baboon amylase locus, we utilized closely related
1005 Papionini and Old World monkey species that still retain the ancestral Catarrhini two-copy configuration
1006 (*AMY2B* and *AMY1'*). Continuous contigs spanning the locus were available for gelada (*Theropithecus*
1007 *gelada*), drill (*Mandrillus leucophaeus*), sooty mangabey (*Cercocebus atys*), and crested macaque
1008 (*Macaca nigra*). We extracted the orthologous intervals (*RNPC3* 5' flank to *AMY1C* 3' flank; see above)
1009 and aligned each of them against the olive baboon amylase locus with NUCmer (v 3.1; default settings),
1010 visualising the results with mummerplot and extracting the corresponding coordinates. Between these
1011 species, uninterrupted collinearity was observed across the junctions that define the "first segment" and
1012 "fourth segment" as defined above. The "second segment" and "third segment" were then mapped
1013 separately against the ancestral two-copy configuration, allowing us to delineate the exact breakpoints
1014 introduced by the first and second non-allelic homologous recombination events. The sequence that is
1015 contiguous in olive baboons but split in the two-copy genomes corresponds to the novel duplications.
1016 These alignments confirmed that no additional rearrangements in the olive baboon amylase locus and
1017 that the derived *AMYp1* and *AMYp2* blocks are absent from the ancestral Catarrhini two-copy
1018 configuration, validating the BISER-based segmentation. We additionally aligned the olive baboon locus
1019 to that of its sister species Guinea baboon (*Papio papio*). Although the Guinea baboon assembly contains
1020 six tandem amylase copies, in addition to the ancestral *AMY2B* and *AMY1'*, all of these extra copies are
1021 identical in nucleotide sequence, strongly suggesting an assembly artifact caused by assembly
1022 over-expansion rather than true structural variation in the locus (see **Methods**). Because Guinea
1023 baboons retain only the ancestral two-copy configuration, both *AMYp1* and *AMYp2* appear to be olive
1024 baboon-specific. Consequently, the two NAHR events that generated these duplications must have
1025 occurred after the complete split from Guinea baboons, within the ~1.85 MYA window inferred for
1026 Papionini divergence

1027 To investigate the macaque-specific duplication, we applied the same segmentation and alignment
1028 workflow to the rhesus macaque locus. Self-alignments (dotplots) and pairwise NUCmer comparisons
1029 partitioned the region into three paralogue segments, each corresponding to a single duplicon: segment 1
1030 harbouring *AMY2B*, segment 2 carrying *AMYm* (the novel copy), and segment 3 harbouring *AMY1'*.
1031 Pairwise comparisons showed that *AMYm* shares extensive 5' homology with *AMY1'* and 3' homology
1032 with *AMY2B*, with crossover points falling within the intergenic sequence upstream the novel gene. The
1033 preserved orientation and high identity of the flanks, together with these shared blocks, implicate NAHR
1034 as the duplication mechanism (**Figure S3**). In parallel, BISER (v. 1.4) identified 26 unique duplicons
1035 within the macaque locus (**Table S13**), which we treat as the smallest repeat units; by contrast, the three
1036 segments defined above represent the largest repeated modules spanning each paralogue block across
1037 the locus.

1038 We then aligned orthologous contiguous amylase locus from multiple *Macaca* species representing the
1039 three major clades (*fascicularis*, *sinica*, *silenus*). *AMYm* is present in the *fascicularis* and *sinica* groups
1040 but absent from the *silenus* group, placing the duplication after the *silenus* split and before diversification
1041 of the *fascicularis* and *sinica* clade, i.e. approximately ~4.5-5 MYA. All *AMYm*-harbouring loci were
1042 recovered on single, gap-free contigs with identical breakpoints across species and no additional
1043 rearrangements, supporting a single NAHR event that generated the three-segment architecture in
1044 macaques (**Figure S4**).

1045 **Targeted Local Reassembly of the Guinea Baboon Amylase Locus**

1046 The original Guinea baboon (*Papio papio*) genome assembly (GCA_028858775.2) contained six identical
1047 tandem amylase copies in addition to *AMY2B* and *AMY1'*. To evaluate whether this expansion reflects
1048 true structural variation or an assembly artifact, we adopted a targeted local reassembly approach. We

1049 extracted an extended bait region spanning from 150 kb upstream of the *RNPC3* 5' flank to 150 kb
1050 downstream of the *AMY1C* 3' flank, providing sufficient unique flanking context for unambiguous read
1051 recruitment. The full PacBio Revio HiFi read dataset (~97 Gb) was mapped to this bait region using
1052 minimap2 (v. 2.29)⁶³, retaining only primary alignments with mapping quality ≥ 20 and minimum aligned
1053 length ≥ 7 kb. Coverage was assessed in 20 kb windows across the extracted locus. A sharp drop in read
1054 depth was observed precisely over the region corresponding to the six tandem copies, inconsistent with a
1055 genuine six-copy expansion and instead indicative of assembly over-expansion.

1056 To confirm this interpretation, the filtered reads were extracted and reassembled locally using both Flye
1057 (v. 2.9.6)⁶⁴ and hifiasm (v. 0.25.0)⁶⁵. Across mapping quality thresholds from MQ 10 to MQ 30 and
1058 minimum aligned lengths from 6 kb to 10 kb, hifiasm consistently produced a single continuous contig
1059 spanning the entire locus, containing only two amylase genes (*AMY2B* and *AMY1'*) corresponding to the
1060 ancestral Catarrhini configuration. Flye occasionally fragmented the locus at tandem repeat boundaries
1061 but never supported more than two gene copies. As an additional control, the analysis was repeated
1062 using the gelada (*Theropithecus gelada*) amylase locus as an independent bait reference; uniform
1063 coverage and a two-copy local assembly were again recovered. Based on the locally reassembled locus,
1064 transposable element annotations for the Guinea baboon amylase region were recalculated using
1065 RepeatMasker with the same parameters applied to all other species, and all downstream analyses were
1066 updated accordingly to reflect the TE composition of the reassembled locus.

1067 **Transcriptomic Data Generation and Processing and Differential Expression Analysis**

1068 Biopsies (parotid, submandibular, sublingual, pancreas and liver) from five olive baboons (*Papio anubis*)
1069 and six rhesus macaques (*Macaca mulatta*) were flash-frozen in liquid nitrogen immediately after
1070 collection and stored at -80°C . These samples were collected at Texas Biomedical Institute by
1071 veterinarians from 7-18 year old animals right after planned euthanization for health reasons. Samples
1072 are kept at -20°C . Library preparation and standard Illumina HiSeq RNA sequencing (2x150bp)
1073 experiments were carried out following standard operating procedures by GENEWIZ/Azenta. Publicly
1074 available adult human salivary gland dataset⁴⁹ was downloaded as FASTQ files and processed similarly
1075 to the olive baboon and rhesus macaque RNAseq samples from the quality-control step onwards. The
1076 pancreas and liver expression datasets were obtained from GTEx v8 and normalized as detailed below.

1077

1078 Quality assessment of raw reads was performed with FastQC (v. 0.12.1)⁶⁶ and summarised with MultiQC
1079 (v. 1.25)⁶⁷. Adapter sequences, low quality bases at the 3' and 5' end of the read, and reads shorter than
1080 36 bp were removed using Cutadapt (v. 3.5)⁶⁸. To investigate tissue-specific expression, trimmed reads
1081 were aligned to the olive baboon (*Papio anubis* Annotation Release 104) and rhesus macaque (*Macaca*
1082 *mulatta* Annotation Release 103) transcriptomes, and transcript abundances were quantified with Kallisto
1083 (v. 0.51.1)⁵⁰. Transcript-level transcripts per million (TPMs) were summarized to the gene level. The
1084 gene-level quantifications were further normalized using the "lengthScaledTPM" method, for consistency
1085 in downstream analyses. These data were then filtered to include only high-confidence gene expression
1086 estimates (minimum of ten reads across all samples). The normalized gene expression profiles were
1087 used for differential expression analysis and visualization. For differential expression analysis, we utilized
1088 DESeq2 (v. 1.46.0)⁶⁹ to compare gene expression profiles between tissues within species. We applied a
1089 Wald test to detect significant expression differences between groups. Genes with an adjusted p-value
1090 (FDR < 0.05) were considered significantly differentially expressed.

1091

1092 Human salivary gland reads⁴⁹ were quantified with Kallisto against the NCBI *Homo sapiens* Annotation
1093 Release 110 transcriptome. Pancreas and liver read-count tables were obtained from GTEx v8. Before
1094 combining these data with the Kallisto-derived human salivary gland quantifications, we removed
1095 Ensembl version suffixes from every gene ID, retained only genes present in all samples and rounded
1096 Kallisto's fractional estimates to integers so that every entry represented a raw count compatible with
1097 DESeq2. A DESeq2 object was then created using tissue as the design factor, and size factors were
1098 estimated with the default median-of-ratios procedure to generate library-normalised counts.
1099 Variance-stabilising transformation was applied for PCA-based quality control (GTEx pancreas, GTEx
1100 liver or each salivary gland tissue) to confirm that batch effects did not dominate the variance structure.
1101 The resulting size-factor-normalised count matrix was used for all subsequent differential-expression
1102 analyses.

1103 In parallel, we developed a polymorphism-aware pipeline to achieve higher resolution quantification of
1104 gene expression across tissues. This pipeline leverages nucleotide polymorphisms within RNA-seq reads
1105 to partition expression into transcript-specific contributions. We first aligned all human, rhesus macaque,
1106 and olive baboon transcripts, accounting for alternative splicing isoforms, to identify coding-sequence
1107 polymorphisms. For each species, we generated a single consensus FASTA sequence that captured the
1108 shared polymorphisms across the different transcripts within species. We then extracted RNA-seq reads
1109 previously mapped to each genome using HISAT2 (v. 2.2.1)⁷⁰, retaining only those overlapping the
1110 amylase locus, and realigned them to the corresponding species-specific consensus FASTA with BWA (v.
1111 0.7.19-r1273)⁷¹. NCBI annotations were curated for known polymorphic regions in the olive baboon and
1112 rhesus macaque transcripts, and the data were integrated with species-specific annotations to refine the
1113 analysis.

1114

1115 Using Old World Monkeys transcripts as outgroups, we assigned the detected polymorphisms as
1116 ancestral, derived, or lineage-specific. Using custom Python scripts we extracted and quantified reads at
1117 these polymorphic sites, enabling paralog-specific expression estimates (scripts have been deposited to
1118 Zenodo:<https://doi.org/10.5281/zenodo.16809248>). By focusing on transcript- and paralog-level variation,
1119 our approach aimed to reveal tissue-specific expression patterns that would otherwise remain obscured
1120 in conventional bulk RNA-seq data. To ensure robust detection of gene-level differences, we specifically
1121 utilized polymorphisms that distinguish one paralogous gene from another. In the rhesus macaque
1122 genome, for instance, three paralogous genes, *AMY2B*, *AMYm* and *AMY1'*, are annotated in that order.
1123 For each paralogous gene, we selected four polymorphisms: two located toward the proximal (5') end
1124 and two toward the distal (3') end of the coding region. At positions 93 and 111, *AMY2B* carries the
1125 polymorphisms G, T, *AMYm* carries A, G, and *AMY1'* carries A, T (**Figure S14**). At positions 737 and 753,
1126 *AMY2B* carries A, T, *AMYm* carries T, T, and *AMY1'* carries T, C respectively. The polymorphism-aware
1127 pipeline allowed us to unambiguously assign reads to individual paralogs and confirmed the accuracy of
1128 Kallisto's expression estimates. It was used only as an internal validation step and did not contribute to
1129 the differential expression analyses or any other downstream analyses reported here.

1130

1131 Finally, we queried all annotated long-non-coding RNAs located within, or immediately flanking, the
1132 amylase locus in the rhesus macaque and olive baboon genome builds. The Kallisto-derived TPMs for
1133 every such lncRNA were below ten reads in every tissue, indicating negligible expression; they therefore
1134 cannot account for the tissue-specific patterns described in Results.

1135 **ddPCR-Based Copy Number Estimation in Individuals Included in the Transcriptomics Analysis**

1136

1137 To determine whether within-species variation in amylase copy number could confound paralog-specific
1138 expression comparisons, we quantified genomic amylase copy number by droplet digital PCR (ddPCR) in
1139 the same five olive baboons (*Papio anubis*) and six rhesus macaques (*Macaca mulatta*) for which
1140 RNA-seq data were generated from parotid, submandibular, sublingual, pancreas, and liver biopsies (see
1141 RNA-seq sample description above). Genomic DNA was isolated from each individual and analyzed
1142 using ddPCR, targeting a sequence conserved across amylase gene homologs in Old World monkeys.
1143 Copy number was calculated relative to a single-copy reference locus used in our previous work¹⁷ (**Table**
1144 **S14**). The assay used the following oligonucleotides: forward primer (5'-3')
1145 GAGCACTTGCTTTGTGGATAA; reverse primer (5'-3') TCCAGAAAGGTAAGAATAGAGG; and
1146 hydrolysis probe (5'-3') CCATGACAATCAACGAGGACATGGG. Because the target region is shared
1147 among paralogs, this approach provides an aggregate estimate of total amylase copy number and does
1148 not resolve paralog-specific contributions to copy number variation.

1149 **Positive Selection Analyses**

1150 Coding sequences (CDS) from the two Old World monkey species analysed here (olive baboons and
1151 rhesus macaques) together with orthologous CDS from all extant Great Ape species (human,
1152 chimpanzee, bonobo, Sumatran and Bornean orangutan) were aligned at the amino-acid level with
1153 MAFFT (v. 7.515)⁷², back-translated to codons with PAL2NAL (v. 14)⁷³, and manually trimmed to preserve
1154 reading frame. We assessed the sequences for internal stop codons and removed any truncated CDS. A
1155 maximum-likelihood gene tree was inferred with IQ-TREE (v 2.4.0)⁷⁴ using the MG+F3X4+R2 model and

1156 was input to all HyPhy analyses. All selection analysis results are provided in Supplementary **Tables**
1157 **S15-S18**.

1158 Branch-level tests (aBSREL)

1159 To identify entire lineages that experienced bursts of adaptive change we ran the adaptive Branch-Site
1160 Random-Effects Likelihood (aBSREL) model in HyPhy (v. 2.5.48)⁷⁵ with default settings. aBSREL
1161 compares for every branch, a null model in which all sites evolve neutrally or under purifying selection (ω
1162 ≤ 1) to an alternative model that allows a proportion of sites on that branch to have $\omega > 1$. Likelihood-ratio
1163 tests are corrected for multiple comparisons with the built-in Holm procedure. Foreground branches were
1164 not pre-specified; instead HyPhy evaluates each branch in turn. The aBSREL p-values were corrected for
1165 multiple testing with HyPhy's built-in Holm-Bonferroni procedure (**Table S15**).

1166 Site-level episodic diversifying tests (MEME)

1167 Codon-specific episodic diversifying selection was evaluated with the Mixed-Effects Model of Evolution
1168 (MEME) in HyPhy⁷⁶. MEME allows the selective regime (ω) at a site to vary among branches, thus
1169 detecting positive selection that operates only on a subset of lineages. We used the default significance
1170 cut-off of $p \leq 0.10$ after applying the False Discovery Rate (FDR) correction using the Benjamini-Hochberg
1171 procedure (**Table S16**).

1172 Site-level pervasive tests (FUBAR)

1173 Pervasive, site-wide selection was assessed with the Fast, Unconstrained Bayesian AppRoximation
1174 (FUBAR) implemented in HyPhy⁷⁷, which estimates posterior probabilities for $\omega > 1$ under a Bayesian
1175 framework that assumes a constant selection regime across the tree. Analyses were run for 5 million
1176 MCMC iterations (burn-in=1 million); codons with posterior probability larger than 0.90 were considered
1177 positively selected (**Table S17**).

1178 Relaxation or intensification of selection (RELAX)

1179 To test variation in selective pressures across specific lineages, we applied RELAX⁷⁸. The baboon *AMYp2*
1180 branch was assigned as foreground and all other branches as background. RELAX fits two models
1181 differing by a scaling parameter K that inflates ($K > 1$) or deflates ($K < 1$) the background ω distribution;
1182 significance is assessed by a likelihood-ratio test (**Table S18**).

1183 **Structural Modeling and Functional Domain Prediction**

1184 The coding sequence for *AMYp2* was modelled *in silico* using AlphaFold2 (v. 2.3.1)⁷⁹ via the ColabFold
1185 implementation with the "monomer_ptm" preset and default recycling. The top-ranked model by pLDDT
1186 was retained; predicted-aligned-error (PAE) matrices were inspected to verify global fold confidence.
1187 Annotated PDB files were generated in Chimera (v. 1.19) and used for all subsequent structure-based
1188 alignments.

1189 Active, catalytic and calcium-binding sites were identified using the NCBI Conserved Domain Database
1190 (CDD, accessed March 2025). Calcium and chloride binding sites were further cross-referenced with the
1191 identified binding sites by Ramasubbu et al.^{80,81}. Glycosylation candidates were predicted with NetNGlyc
1192 1.0⁸² and cross-referenced to the proposed glycosylation sites reported by Kamitaki et al.⁴⁵.

1193 **Regulatory Motif Analysis**

1194 Promoter sequences for the amylase paralogs in humans and rhesus macaques were defined as the
1195 170-bp window spanning 100 bp upstream to 70 bp downstream of the experimentally supported TSS
1196 recorded in Eukaryotic Promoter Database (EPD)⁸³. When an EPD entry was unavailable, we took the
1197 RefSeq transcription-start site (TSS) from the corresponding annotation (NCBI *Papio anubis* Release
1198 104, *Macaca mulatta* Release 103, and *Homo sapiens* Release 110) and defined the promoter as the
1199 region 100 bp upstream to 70 bp downstream of that TSS. We retrieved these windows for every amylase
1200 paralog in humans, rhesus macaques, and olive baboons. In parallel, we analyzed the 50 most highly
1201 expressed salivary gland genes in each species (ranked by TPM) and annotated their promoter
1202 sequences. Each promoter was scanned with MEME (MEME-suite v. 5.5.7)⁸⁴ for *de novo* motif discovery
1203 to identify novel, enriched sequence motifs.

1204 The identified motifs were then annotated with Tomtom (MEME-suite v. 5.5.8)⁸⁵ against the JASPAR 2024
1205 CORE non-redundant vertebrate library⁸⁶, retaining hits with $P < 10^{-4}$ to identify potential transcription
1206 factors (TFs). To assess specificity, we conducted parallel analyses using promoter sequences from 50
1207 randomly selected genes per species, which did not show significant motif enrichment for the motifs
1208 identified with the salivary gland dataset, supporting the specificity of our results. The resulting TF list was
1209 cross-referenced with salivary- and pancreas-specific transcription factors with enriched expression from
1210 the FANTOM5 database ($P < 0.05$ and $\log_{10}(\text{relative expression over median}) > 1.3$)⁸⁷ and the salivary gland
1211 TF catalogue from Michael et al.⁵¹. TFs present in either salivary gland set and absent from the
1212 pancreatic set were labelled salivary-gland-biased, while the converse defined pancreatic-biased TFs,
1213 with the remainder classed as core. The amylase-paralog promoter windows were then scanned with
1214 FIMO⁸⁸ against the JASPAR 2024 CORE library to identify TFBS within each promoter; only hits with
1215 $P < 10^{-4}$ were retained. These TFBS were subsequently grouped into core, pancreatic-biased, and
1216 salivary-gland-biased categories based on the above TF assignments.

1217 **ADDITIONAL RESOURCES**

1218 Analysis scripts, pipelines and input data: <https://doi.org/10.5281/zenodo.16809248> and
1219 <https://doi.org/10.5281/zenodo.18689074>