1

2

3

# Recurrent rearrangements of human amylase genes create multiple independent CNV series

Nzar A.A. Shwan[1,3*], Sandra Louzada[2*], Fengtang Yang[2] and John A.L. Armour[1]

[* = co-first author]

1 School of Life Sciences, University of Nottingham, Medical School, Queen's Medical Centre, Nottingham NG7 2UH, UK

2 Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, UK

3 Scientific Research Centre, University of Salahaddin, Erbil, Kurdistan, Iraq

Corresponding author: John Armour (john.armour@nottingham.ac.uk)

*Key words: genomic mutation; adaptation; genomic instability; CNV*

19  **Abstract**

20  The human amylase gene cluster includes the human salivary (*AMY1,* MIM#

21  104700) and pancreatic amylase genes (*AMY2A*, MIM# 104650 and *AMY2B,* MIM#

22  104660), and is a highly variable and dynamic region of the genome. Copy number

23  variation of *AMY1* has been implicated in human dietary adaptation, and in

24  population association with obesity, but neither of these findings has been

25  independently replicated. Despite these functional implications, the structural

26  genomic basis of copy number variation (CNV) has only been defined in detail very

27  recently. In this work we use high-resolution analysis of copy number, and analysis

28  of segregation in trios, to define new, independent allelic series of amylase CNVs in

29  sub-Saharan Africans, including a series of higher-order expansions of a unit

30  consisting of one copy each of *AMY1*, *AMY2A* and *AMY2B.* We use fibre-FISH

31  (fluorescence *in situ* hybridization) to define unexpected complexity in the

32  accompanying rearrangements. These findings demonstrate recurrent involvement

33  of the amylase gene region in genomic instability, involving at least five independent

34  rearrangements of the pancreatic amylase genes (*AMY2A* and *AMY2B*). Structural

35  features shared by fundamentally distinct lineages strongly suggest that the common

36  ancestral state for the human amylase cluster contained more than one, and

37  probably three, copies of *AMY1*.

38

39

**Introduction**

40

41   The adoption of agriculture was one of the most radical and pervasive innovations

42   among the many changes introduced by humans to their own environments. In

43   addition to a capacity to support higher population densities, agricultural food

44   production has led to a shift in dietary composition, including increases in dietary

45   starch as the result of reliance on starch-rich staples. Starch is initially digested by

46   the enzyme amylase, present in humans in two tissue-specific isoenzymes: salivary

47   amylase, encoded by the gene *AMY1*, and pancreatic amylase, encoded by *AMY2A*

48   and *AMY2B*. These amylase genes are all found in a cluster on human chromosome

49   1, and early observations on pedigree segregation of protein electrophoretic variants

50   demonstrated common and extensive copy number variation (CNV) in the salivary

51   amylase gene *AMY1* [Pronk and Frants, 1979; Pronk et al., 1982]. These

52   observations, coupled with detailed mapping of cloned genomic sequences, showed

53   that there were common haplotypes containing odd numbers of *AMY1* genes,

54   differing by pairs of genes in inverted orientation [Bank et al., 1992; Groot et al.,

55   1989; Groot et al., 1991; Groot et al., 1990]. More recently, higher-resolution studies

56   of the variation have demonstrated that most humans have an even number of

57   *AMY1* copies, as predicted by the predominance of haplotypes containing odd

58   numbers, with an overall copy number range of 2 to 18 copies per individual

59   [Carpenter et al., 2015; Usher et al., 2015].

60   Primarily because of its early discovery and extensive range, most attention on

61   amylase CNVs has focussed on the salivary amylase gene *AMY1*, but there have

62   been reports of CNVs involving the *AMY2* genes [Conrad et al., 2010; Cooke Bailey

63   et al., 2013; Groot et al., 1991; Sudmant et al., 2010]. Integration of information from

64   read-depth analysis, segregation and direct typing of copy number demonstrated

65   haplotypes harbouring even numbers of *AMY1* in conjunction with CNVs of the

66   pancreatic amylase genes *AMY2A* and *AMY2B*. There are two common CNVs of

67   *AMY2* genes in European populations – one carrying a deletion of the *AMY2A* gene,

68   the other a duplication of both *AMY2A* and *AMY2B* [Carpenter et al., 2015; Usher et

69   al., 2015]. Those investigations also implied that there were other rearrangements of

70   the locus that could not be accounted for by the allelic series common in Europe.

71   The extensive variation in *AMY1* copy number has prompted studies exploring its

72   functional significance, including the observation that populations with starch-rich

73   diets appear to have significantly higher average *AMY1* copy number than

74   populations with lower starch intake [Perry et al., 2007]. The implication that copy

75   number expansion of *AMY1* represents an adaptation to dietary shifts following the

76   adoption of agriculture fits with the observation that the gene is found as a single

77   copy in chimpanzees [Perry et al., 2006], and in the genomes of archaic hominins

78   [Lazaridis et al., 2014; Olalde et al., 2014]. More recently, the observation of a

79   significant correlation between low *AMY1* copy number and higher body mass index

80   (BMI) suggested that the CNV had considerable ongoing functional importance in

81   modern humans [Falchi et al., 2014]. Although further studies have supported the

82   association [Mejía-Benítez et al., 2015], doubt was cast on the medical importance of

83   the association by the failure of a rigorous and well-powered study to reproduce the

84   observation [Usher et al., 2015]. Most recently, a carefully calibrated study of *AMY1*

85   copy number in East Asian samples also failed to demonstrate any association with

86   BMI [Yong et al., 2016].

87   In this work we set out to understand more thoroughly the range of common genomic

88   variation in amylase copy number found in humans, and in particular to define the

89   potential range of CNVs of *AMY2* genes. We combine high-resolution DNA typing,

90 fibre-FISH and SNP analysis to show that independent rearrangements of the *AMY2*

91 genes have arisen on at least five occasions, and can include haplotypes containing

92 up to 5 copies each of *AMY2A* and *AMY2B*. Although we cannot exclude neutral

93 mutation processes at high frequency in this highly repetitive and unstable region,

94 recurrent and human-specific rearrangements suggest the likelihood of adaptive

95 value for these variants.

96

97   **Materials and Methods**

98   *Amylase copy number determination*

99   Previously published methods were used to measure relative representation of

100   *AMY1*-coupled microsatellite alleles and ratios of *AMY2A*:*AMY2B* copy numbers

101   [Carpenter et al., 2015]. *AMY1* copy number was measured by modified PRT

102   approaches, in which distinctive sequence variants from the two terminal

103   (centromeric) copies of *AMY1* ("*AMY1*C") were used as reference loci. Two

104   fluorescent PCRs in a total volume of 10μl were done per sample, each using three

105   primers at 1μM and 10ng genomic DNA in the buffer described [Carpenter et al.,

106   2015], and switching the activities of primers using cycling conditions. The first PRT

107   uses primers AMY1CF, HEX-AMY1CR and nested forward primer NF2, and the

108   second contains AMY1CF, FAM-labelled AMY1CRB2 and nested forward primer

109   NF5 (see Table 1).

110

111   Reactions started with 15 cycles of 95°C 30s/ 61°C 30s/ 65°C 2 minutes, during

112   which AMY1CF and AMY1CR/RB2 anneal stably to make products specific to

113   *AMY1*. The cycles then switched to 95°C 30s/ 54°C 30s/ 65°C 1 minute, for 14

114   (AMY1CR+NF2) or 13 (AMY1CRB2/NF5) cycles, before final extension at 72°C for

115   50 minutes; at the lower annealing temperature in the second phase the nested

116   primers NF2 and NF5 anneal stably to make shorter products that are more readily

117   resolved. PCR products were quantified after separation by capillary electrophoresis

118   on an ABI3130xl Genetic Analyser 36 cm capillary, running the products from the

119   two reactions in the same capillary. Before electrophoresis 2μl from reactions with

120   AMY1CR/NF2 and 0.8μl from reactions with AMY1CRB2/NF5 were mixed in 10μl

121   HiDi formamide containing 0.125μl ROX-500 markers. These samples were

122   denatured at 96°C for 3 minutes before electrophoresis using POP- 7 polymer and

123   an injection time of 30s at 1kV. GeneMapper software (Applied Biosystems) was

124   used to extract peak area data.

125

126   In nearly all samples AMY1CR/NF2 amplify 436bp products from the two *AMY1C*

127   copies and 427bp products from all other (*AMY1A/1B*) copies.

128   AMY1CRB2/NF5 amplify 357bp products from typical copies of *AMY1C* and 344bp

129   products from *AMY1A/1B*; a distinctive alternative product of 347bp is amplified from

130   the variant *AMY1C* haplotype. Ratios of *AMY1A1B* to *AMY1C* can be used to

131   deduce *AMY1* copy number, assuming that there are two copies of *AMY1C*,

132   calibrating the data with integer clusters defined using k-means clustering. Further

133   details and representative data can be found in the Supplementary Material and

134   Supp. Figures S1-S5.

135

136   The assay for the *AMY2A/2B* duplication junction fragment [Carpenter et al., 2015]

137   was modified to allow quantitative readout after capillary electrophoresis of

138   fluorescent PCR products. PCRs of 10µl used 10ng genomic DNA in the buffer

139   described [Carpenter et al., 2015], with final concentrations of 1µM of each of three

140   primers AMY2B2D, FAM-AMY2B2R and AMY2B2F (Table 1).

141

142   PCRs used an initial denaturation stage of 95°C for 5 minutes, followed by 22 cycles

143   of 95°C 30s/ 60°C 30s/ 65°C 1 minute, and final extension at 72°C for 50 minutes.

144   Products of 192bp between AMY2BF and AMY2BR are made from all samples, and

145   if it is present the duplication junction sequence produces a 176bp product between

146   AMY2BD and AMY2BR. Before electrophoresis 1µl from PCRs was mixed with 10µl

147     HiDi formamide containing 0.125µl ROX-500 markers, and denatured and separated

148     by capillary electrophoresis as above.

149

150     Copy number ratios for *AMY1* relative to *AMY2* (*AMY2A*+*AMY2B*) were determined

151     by a PRT exploiting a consistent 4bp length difference in the paralogous products

152     from just upstream of exon 4, using primers HEX- AMY1_2F  and AMY1_2R (Table

153     1). The ratios of products from *AMY1* (169bp) to *AMY2A* + *AMY2B* (173bp) were

154     used to infer the ratio of genomic copy numbers.

155

156

157     *Fibre-FISH methods*

158     The probes and general methods for fibre-FISH are given in detail in [Gribble et al.,

159     2013] and [Carpenter et al., 2015]. In summary, DNA fibres were prepared from

160     agarose-embedded cells by molecular combing (Genomic Vision), and probes were

161     derived from one PCR product from the *AMY1* gene [Perry et al., 2007], and one

162     each from the regions upstream of *AMY2A* and *AMY2B* [Carpenter et al., 2015].

163

164     *Haplotype definition and database records*

165     In the Leiden Open Variation Database format (LOVD, http://www.lovd.nl [Fokkema

166     et al., 2011]), information defining structural allelic variants involving the *AMY1*,

167     *AMY2A* and *AMY2B* genes is collected under the *AMY2B* locus-specific database

168     (http://www.LOVD.nl/AMY2B). In this work we have given the *AMY2B* locus-specific

169     database ID of each new or known haplotype – for example, the

170     $(AMY1)_3(AMY2A)_1(AMY2B)_1$ haplotype found in the human reference assembly hg19

171     has the ID AMY2B_011111.

**Results**

Most haplotypes of the human amylase genes include one copy each of the *AMY2B* and *AMY2A* genes and an odd number of copies of *AMY1*. The differing *AMY1* copy numbers arise from variation in the numbers of a 95kb cassette including two copies of *AMY1* and one copy of the truncated *AMY2A* pseudogene "*AMYP1*" [Carpenter et al., 2015; Usher et al., 2015]. The arrangement of the sequence in the hg19 human reference assembly conforms to this pattern (locus-specific database (http://www.LOVD.nl/AMY2B) ID AMY2B_011111), with three copies of *AMY1*, and is illustrated in the upper panel of Figure 1. A common haplotype pattern not conforming to this structure has been described in recent work [Carpenter et al., 2015; Usher et al., 2015] (database ID AMY2B_022101); in this, *AMY2B*, *AMY2A* and one copy of *AMY1* are duplicated (via non-homologous rearrangement), creating a unique junction (shown as "J" in Figure 1) that can form the basis of a PCR assay for the structure [Carpenter et al., 2015]. Fibre-FISH confirmation of this structure for European sample GM12239 is shown in Supp. Figure S6. These *AMY2A2B* duplications are characteristically associated with haplotypes containing even numbers of *AMY1* (usually 4), and are common in European and African populations, but less so in East Asians [Carpenter et al., 2015; Usher et al., 2015]; in the nomenclature of Usher *et al.* [Usher et al., 2015], two examples are designated AH2B2 (AMY2B_022200) and AH4B2 (AMY2B_022211).

*Higher-order expansions of pancreatic amylase genes*

To understand the full scope of variation in human amylase genes, we aimed first to define the composition and structures of alleles containing more than two copies of each of the pancreatic amylase genes *AMY2A* and *AMY2B*. The gene content of

9

196  haplotypes in Yoruban (YRI) trios from the HapMap phase 1 were determined first by

197  measuring the gene copy numbers of *AMY1*, *AMY2A* and *AMY2B*, followed by

198  analysis of segregation of *AMY1*-coupled microsatellite alleles [Carpenter et al.,

199  2015], *AMY1:AMY2* ratios and *AMY2A:2B* ratios in trios. For most parental samples

200  our direct measurements (Supplementary Dataset) were corroborated by read-depth

201  measures [Carpenter et al., 2015]. For application in this work we developed a new

202  PRT method to measure *AMY1* copy number based on the ratio between distinctive

203  sequences at the centromeric (*AMY1*C) copy and the internal (*AMY1*A/B) copies; in

204  practice, we found that this measure combined high levels of accuracy with the

205  convenience of assigning most samples to integer classes with no more than two

206  PCRs (Figure 2a and Supplementary Material). In parallel, we modified our assay for

207  the junction sequence specific to the *AMY2A2B* duplication allele, to allow

208  quantification of that sequence relative to the diploid genome (Figure 2b).

209  In most cases (see Figure 3 and Table 2) measurement of copy numbers and ratios

210  in Yoruban trios allowed deduction of the likely haplotype composition. The copy

211  number data were consistent with analyses based on read-depth from the 1000

212  Genomes Project [Carpenter et al., 2015; Usher et al., 2015], and demonstrate that

213  there are distinctive haplotypes associated with higher-order amplifications of

214  *AMY2A* and *AMY2B*, including triplication, quadruplication and quintuplication

215  (AMY2B_033201/044301/055401); in nearly all cases, alleles carrying higher-order

216  expansions of *AMY2A* and *AMY2B* are predicted to carry equal numbers of *AMY1*,

217  *AMY2A* and *AMY2B* genes, so that (for example) the untransmitted maternal

218  quintuplication allele in family Y056 has the composition

219  $(AMY1)_5(AMY2A)_5(AMY2B)_5$ (AMY2B_055401, Table 2). Quantification of product

220  ratios showed that n-fold expansions of (*AMY2B-AMY2A-AMY1*) contained (n-1)

221     copies of the junction sequence found in the *AMY2A+2B* duplication allele series

222     ([Carpenter et al., 2015]; AMY2B_022101 above, or AMY2B_022200 and

223     AMY2B_022211, equivalent to alleles AH2B2 and AH4B2 in [Usher et al., 2015]).

224     This observation suggested that there is a new allelic series based on higher

225     expansion of the repeat unit formed in the *AMY2A+2B* duplication allele, with the

226     known junction sequence separating adjacent copies of an (*AMY2B-AMY2A-AMY1*)

227     repeat unit.


228     We applied fibre-FISH to define the physical structure of the haplotypes we had

229     defined on the basis of gene content; our previous observations demonstrated

230     [Carpenter et al., 2015] that although the high level of sequence similarity between

231     amylase gene sequences leads to cross-hybridization, especially between *AMY1*

232     and *AMY2A*, it is nevertheless possible to distinguish the *AMY1* and *AMY2A* genes

233     on the basis of hybridization patterns (Figure 3, top). By contrast, the sequence

234     upstream of *AMY2B* is sufficiently distinct to give locus-specific hybridization. Fibre-

235     FISH analysis of expanded alleles verified the prediction of a repeat unit containing

236     one copy of each gene, but also showed that in all cases the first (telomeric) unit

237     contained an inversion, to give the gene order (*AMY2B-AMY1-AMY2A*)-(*AMY2B-*

238     *AMY2A-AMY1*)$_{(n-1)}$. This observation suggested the detailed structure for the

239     triplication allele (AMY2B_033201) in family Y060 (Table 3) shown in Figure 3. The

240     inversion of the first telomeric unit is also seen in fibre-FISH analysis of

241     quadruplication and quintuplication alleles (Supp. Figures S7 and S8), but escaped

242     detection by optical mapping [Usher et al., 2015]. We used long PCR and Sanger

243     sequencing (Supplementary Material) to amplify a 9.8kb product across this

244     inversion in the quintuplication (AMY2B_055401) carrier NA19159 (GenBank

245     KX394682). This sequence verified the orientations shown in Figure 3, but

246    demonstrated no further rearrangements or sequence variants unique to this

247    structure.

248    Alleles containing higher-order (n ≥ 3) expansions of *AMY2A* and *AMY2B* were

249    examined by fibre-FISH (4 examples), segregation (4 examples) and analysis of

250    1000 Genomes Project read-depth (12 examples of AFR individuals with more than

251    3 copies of both *AMY2A* and *AMY2B*); these alleles appeared to be coherent for

252    general structure, gene content and SNP associations. All examples

253    (AMY2B_033201/044301/055401) of higher-order amplifications of the unit (*AMY2B-*

254    *AMY2A-AMY1*) were associated (D' = 1) in African populations with a common

255    haplotype tagged by (for example) the derived allele rs12075086T, the same

256    haplotype associated with simple duplication (AMY2B_022101) of *AMY2A* and

257    *AMY2B* in worldwide populations [Carpenter et al., 2015; Usher et al., 2015].

258    Consistent with this conclusion of a single origin for all alleles containing

259    amplifications of both *AMY2A* and *AMY2B*, all contained the same junction sequence

260    (see Methods), and the deduced *AMY1* microsatellite allele content of expanded

261    alleles resembled each other, and those of the duplication allele AMY2B_022101,

262    with a predominance of microsatellite alleles yielding PCR products of 269bp

263    (Supplementary Dataset and [Carpenter et al., 2015]).

264

265    *Duplication of AMY2A*

266    Our previous work and that of others demonstrated individuals (and therefore

267    haplotypes) with higher numbers of *AMY2A* than *AMY2B* [Carpenter et al., 2015;

268    Usher et al., 2015]. Such individuals are more frequent in African populations than

269    others; for example, in our read-depth analysis of 1000 genomes samples, 13.6% of

270 African samples had more copies of *AMY2A* than *AMY2B*, compared with 2.51% of

271 Asians and 0.55% of Europeans [Carpenter et al., 2015]. Segregation analysis in

272 African (YRI) trios confirmed the prediction that the corresponding haplotypes carried

273 a duplication of *AMY2A* unaccompanied by duplication of *AMY2B* (Table 3). As

274 predicted for independently-arising duplications, they were not associated with the

275 specific junction fragment characteristic of the *AMY2A+2B* duplication haplotype.

276 Analysis of SNP associations in these individuals suggest that most examples of

277 *AMY2A*-only duplications were found on a single haplotype background, but there

278 was also evidence of heterogeneity, with (for example) NA19119, who has both

279 haplotypes with an *AMY2A*-only duplication (AMY2B_012341 and AMY2B_012211)

280 on two different SNP haplotypes. We undertook fibre-FISH analysis in family trio

281 Y060 (Table 3), in which both haplotypes in the father NA19119 were predicted to

282 have 2 copies of *AMY2A* and a single copy of *AMY2B* (Figure 3 and Supp. Figure

283 S9).

284 The haplotypes characterised have structures that do not require the formation of

285 new junctions, and can be created by new juxtapositions of sequences present in the

286 reference haplotype AMY2B_011111. However, the structural differences indicate

287 that these two haplotypes (AMY2B_012341 and AMY2B_012211) arose

288 independently of one another,, and that the gene content feature common to these

289 two haplotypes, amplification of *AMY2A* without amplification of *AMY2B*, appears

290 coincidental rather than as the result of common ancestry. In particular, the shorter

291 $(AMY1)_4(AMY2A)_2(AMY2B)_1$ haplotype (AMY2B_012211) has the duplicated copy of

292 *AMY2A* in inverted orientation (Supp. Figure S9). Comparison of amylase copy

293 number with flanking SNP haplotypes suggests that this $(AMY1)_4(AMY2A)_2(AMY2B)_1$

294 haplotype (AMY2B_012211) is the commonest type of $(AMY2A)_2(AMY2B)_1$ structure,

295    whereas we found no evidence for other alleles corresponding to the longer

296    (*AMY1*)$_8$(*AMY2A*)$_2$(*AMY2B*)$_1$ haplotype (AMY2B_012341) in NA19119.

297

298    *A new AMY1:AMY2A junction*

299    Our analysis of copy number segregation in family Y072 consistently indicated

300    ambiguous copy number of *AMY1* in the mother NA19152 and her child NA19154

301    (Table 3); measures of *AMY1* based on the microsatellite (upstream of the *AMY1*

302    gene) indicated 3 copies in the transmitted maternal haplotype, but only 2 copies

303    based on the (downstream) PRT, and intermediate values based on read depth

304    analysis of 1000 Genomes Project reads from the mother NA19152 (estimates of 5.2

305    and 5.55 from [Carpenter et al., 2015] and [Usher et al., 2015] respectively).

306    Segregation also indicated that this haplotype (AMY2B_023201) contained 3 copies

307    of *AMY2A* and 2 copies of *AMY2B* (Table 3). Fibre-FISH analysis confirmed the

308    overall composition of the haplotype, but also demonstrated a hybrid structure with a

309    copy of *AMY2A* and its upstream sequence immediately interrupting one copy of

310    *AMY1* (Supp. Figure S10). Examination of 1000 Genomes Project data from

311    NA19152 showed a single read conforming to this hybrid junction, from which PCR

312    primers were used to demonstrate that the new junction interrupted *AMY1* in exon 4,

313    with 3bp microhomology at the breakpoint (GenBank KX230759).

314

315

**Discussion**

Our work shows that pancreatic amylase (*AMY2A/2B*) genes appear to have

undergone at least five independent rearrangements to create new copy numbers in

humans since the split from chimpanzees. The first is a seamless deletion of *AMY2A*

(AMY2B_010011) common in Europeans and to a lesser extent in Africans, and

generally found on a single SNP background [Carpenter et al., 2015; Usher et al.,

2015]. The second is a duplication of *AMY2A* and *AMY2B* common in Europeans

and Africans (AMY2B_022101/022200/022211), that results from a non-homologous

duplication of an *AMY2B/AMY2A/AMY1* unit, again associated with a common SNP

background [Carpenter et al., 2015; Usher et al., 2015]. Our data in this study show

that this *AMY2A/2B* duplication rearrangement was the starting-point for higher-order

homologous expansions of *AMY2A/2B* found in African populations, as exemplified

by the triplication, quadruplication and quintuplication haplotypes

(AMY2B_033201/044301/055401) we have characterised (Figure 3, Supp. Figures

S7 and S8). There are at least two further lineages

(AMY2B_012211/AMY2B_012341) with independent homologous exchanges

resulting in duplication of *AMY2A* without concomitant duplication of *AMY2B* (Figure

3, Supp. Figure S9), and finally a fifth (non-homologous) rearrangement in which one

copy of *AMY1* is interrupted at exon 4 by a duplication of *AMY2A* (AMY2B_023201,

Supp. Figure S10).

In addition to these rearrangements involving *AMY2A/2B*, allelic series differing by

the 95kb unit containing two repeats of *AMY1* create further overall structural

diversity [Carpenter et al., 2015; Usher et al., 2015]. To summarise the different

mechanisms that have operated in the generation of diversity at this locus in

humans, there have been apparently homologous deletions or duplications of the

15

341   95kb (*AMY1A-AMY1B-AMYP1*) unit, and unequal recombination between

342   homologous repeats 75kb apart is involved in the generation of the *AMY2A* deletion

343   allele. By contrast, the duplication of the 116kb (*AMY2B-AMY2A-AMY1*) unit shows

344   no evidence of being mediated by sequence similarity. Once the duplication is

345   established, however, the generation of higher-order repeats of the (*AMY2B-*

346   *AMY2A-AMY1*) unit could be generated by unequal exchanges between cognate

347   sequences in 116kb repeat sequences. Without complete allele sequences,

348   however, it is difficult to exclude the possibility that additional complexity is involved

349   in some of the apparently simple exchanges between repeats. From a

350   methodological standpoint it is noteworthy that some features of our findings,

351   including the overall structure of the haplotypes, could not be defined using short-

352   read sequencing alone. Long-read capabilities exceeding 10kb would be needed to

353   resolve features, such as the inversion accompanying higher order expansion of

354   *AMY2A* and *AMY2B*, which are clearly demonstrated by fibre-FISH (Figure 3, Supp.

355   Figures S7 and S8), and even then it is unlikely that the overall spatial organisation

356   of the 116kb (*AMY2B-AMY2A-AMY1*) units could be reconstructed unambiguously

357   by primary read assembly, especially if both haplotypes in an individual were of

358   unknown structure.

359   Where genomic rearrangements involve repeated sequences across scales

360   refractory to direct characterisation by sequence assembly of short fragments,

361   longer-range methods such as fibre-FISH or pulsed-field gel electrophoresis can be

362   used to establish haplotype structures. As demonstrated here, fibre-FISH, using

363   combed single-molecular DNA fibres, enabled us to resolve the order, orientation

364   and copy number of amylase family genes on each haplotype unambiguously, even

365   without the need to analyse all three members of each family trio.  However, these

16

366    methods do not provide detailed DNA sequence information. Large-insert (fosmid or

367    BAC) cloning can be used to recover both DNA sequence and information about

368    long-range spatial organisation; it still remains particularly difficult to reconstruct full

369    haplotype sequences when there is population structural allelic variation, as in some

370    disease-associated rearrangements at structurally variable sites, such that in any

371    one sample *both* copies are of unknown structure, for example [Carvalho and Lupski,

372    2008; Yuan et al., 2015], and this study.

373    Including the well-established allelic series differing in copy numbers of *AMY1*

374    (AMY2B_011111/AMY2B_011100/AMY2B_011122, etc.) [Carpenter et al., 2015;

375    Groot et al., 1989; Groot et al., 1990; Usher et al., 2015], it is clear that structural

376    diversity at the human amylase locus has arisen by both homologous and non-

377    homologous events, and has involved rearrangements of both the salivary (*AMY1*)

378    and pancreatic (*AMY2A* and *AMY2B*) amylase genes. The spread of independently-

379    arising rearrangements of the locus can be seen as consistent with the proposal that

380    higher copy-number alleles have been selectively advantageous specifically in

381    recent human history, as suggested by apparently recent human-specific

382    amplification from the single-copy state represented in modern chimpanzees and the

383    genomes of archaic hominins [Lazaridis et al., 2014; Olalde et al., 2014; Prufer et al.,

384    2014].

385    However, it is noteworthy that all the major allelic series of human amylase CNVs

386    defined to date share evidence of the rearrangement that gave rise to the inverted

387    copy of *AMY1* ("*AMY1B*") and the corresponding intergenic region ("18kb" in Figure

388    1), suggesting that the ancestral state for modern humans must have had multiple

389    copies of *AMY1*. The amylase cluster is a region of late-replicating DNA, and is

390    therefore predicted to be prone to frequent rearrangement [Koren et al., 2012; Usher

391  et al., 2015]. Germline mutation to create new copy number alleles cannot be scored

392  simply by observing copy number mismatch in family trios, and first requires enough

393  segregation information to define the parental haplotypes unambiguously; we have

394  nevertheless screened 440 microsatellite haplotype transmissions in three-

395  generation (CEPH) pedigrees without observing any changes in copy number state,

396  suggesting a germline mutation frequency below 0.7% (J.A. and Andrew Cubbon,

397  unpublished work). Given the appearance of similar structures on diverse modern

398  human haplotype backgrounds, the most recent common ancestral state of the locus

399  for all humans is likely to have contained not one copy of each gene, as found in

400  chimpanzees, but instead a sequence similar to the hg19 reference assembly

401  structure $(AMY1)_3(AMY2A)_1(AMY2B)_1$ (AMY2B_011111, equivalent to "H1" of Groot

402  *et al.* [Groot et al., 1989; Groot et al., 1990], or "AH3" of Usher et al. [Usher et al.,

403  2015]). This structure already contained both inverted and tandem-repeated

404  sequences that could predispose to further recurrent rearrangement in the germline,

405  and was itself the result of a non-homologous rearrangement.

406  If an $(AMY1)_3$ allele was the common ancestral structure for all modern humans, the

407  initial amplification to higher gene copy number may have been selectively

408  advantageous before the neolithic, consistent with a recent analysis of sequence

409  data [Inchley et al., 2016]. Nevertheless, whether adaptive or neutral, a preneolithic

410  expansion to higher copy number does not itself preclude subsequent adaptive value

411  for copy number change after the neolithic [Perry et al., 2007].

412  We present no association data relevant to the potential influence of this CNV on

413  obesity, but our results still have implications for the design and interpretation of

414  such studies. Specifically, the expansions of *AMY2* genes we describe here suggest

415  that any influence of amylase gene copy number on body fat is likely to have

416    different genetic architecture in individuals of recent African ancestry. More

417    generally, the extensive structural allelic diversity at the amylase CNV emphasises

418    the extreme difficulty of imputing allelic diversity from SNP data, or of reconstructing

419    structural alleles based on short-read sequence data.

420

421

422

423

424    **Acknowledgements**

429

**References**

Bank RA, Hettema EH, Muijs MA, Pals G, Arwert F, Boomsma DI, Pronk JC. 1992. Variation in gene copy number and polymorphism of the human salivary amylase isoenzyme system in Caucasians. Hum Genet 89(2):213-222.

Carpenter D, Dhar S, Mitchell L, Fu B, Tyson J, Shwan N, Yang F, Thomas MG, Armour JAL. 2015. Obesity, starch digestion and amylase: Association between copy number variants at human salivary (AMY1) and pancreatic (AMY2) amylase genes. Hum Mol Genet 24:3472-3480.

Carvalho CMB, Lupski JR. 2008. Copy number variation at the breakpoint region of isochromosome 17q. Genome Res 18(11):1724-1732.

Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, Aerts J, Andrews TD, Barnes C, Campbell PJ, Fitzgerald T, Hu M et al. 2010. Origins and functional impact of copy number variation in the human genome. Nature 464:704-712.

Cooke Bailey JN, Lu L, Chou JW, Xu J, McWilliams DR, Howard TD, Freedman BI. 2013. The role of copy number variation in African Americans with type 2 diabetes-associated end stage renal disease. J Mol Genet Med 7:61.

Falchi M, El-Sayed Moustafa JS, Takousis P, Pesce F, Bonnefond A, Andersson-Assarsson JC, Sudmant PH, Dorajoo R, Al-Shafai MN, Bottolo L, Ozdemir E, So H-C et al. 2014. Low copy number of the salivary amylase gene predisposes to obesity. Nat Genet 46:492-497.

Fokkema IFAC, Taschner PEM, Schaafsma GCP, Celli J, Laros JFJ, den Dunnen JT. 2011. LOVD v.2.0: the next generation in gene variant databases. Hum Mutat 32(5):557-563.

454 Gribble SM, Wiseman FK, Clayton S, Prigmore E, Langley E, Yang F, Maguire S, Fu

455     B, Rajan D, Sheppard O, Scott C, Hauser H et al. 2013. Massively Parallel

456     Sequencing Reveals the Complex Structure of an Irradiated Human

457     Chromosome on a Mouse Background in the Tc1 Model of Down Syndrome.

458     PLoS ONE 8(4):e60482.

459 Groot PC, Bleeker MJ, Pronk JC, Arwert F, Mager WH, Planta RJ, Eriksson AW,

460     Frants RR. 1989. The Human Alpha-Amylase Multigene Family Consists of

461     Haplotypes with Variable Numbers of Genes. Genomics 5(1):29-42.

462 Groot PC, Mager WH, Frants RR. 1991. Interpretation of polymorphic DNA patterns

463     in the human alpha-amylase multigene family. Genomics 10(3):779-785.

464 Groot PC, Mager WH, Henriquez NV, Pronk JC, Arwert F, Planta RJ, Eriksson AW,

465     Frants RR. 1990. Evolution of the human alpha-amylase multigene family

466     through unequal, homologous, and interchromosomal and intrachromosomal

467     crossovers. Genomics 8(1):97-105.

468 Inchley CE, Larbey CDA, Shwan NAA, Pagani L, Saag L, Antão T, Jacobs G,

469     Hudjashov G, Metspalu E, Mitt M, Eichstaedt CA, Malyarchuk B et al. 2016.

470     Selective sweep on human amylase genes postdates the split with

471     Neanderthals. Scientific Reports 6:37198.

472 Koren A, Polak P, Nemesh J, Michaelson Jacob J, Sebat J, Sunyaev Shamil R,

473     McCarroll Steven A. 2012. Differential Relationship of DNA Replication Timing

474     to Different Forms of Human Mutation and Variation. Am J Hum Genet

475     91(6):1033-1040.

476 Lazaridis I, Patterson N, Mittnik A, Renaud G, Mallick S, Kirsanow K, Sudmant PH,

477     Schraiber JG, Castellano S, Lipson M, Berger B, Economou C et al. 2014.

478        Ancient human genomes suggest three ancestral populations for present-day

479        Europeans. Nature 513(7518):409-413.

480    Mejía-Benítez M, Bonnefond A, Yengo L, Huyvaert M, Dechaume A, Peralta-Romero

481        J, Klünder-Klünder M, García Mena J, El-Sayed Moustafa J, Falchi M, Cruz

482        M, Froguel P. 2015. Beneficial effect of a high number of copies of salivary

483        amylase AMY1 gene on obesity risk in Mexican children. Diabetologia

484        58(2):290-294.

485    Olalde I, Allentoft ME, Sanchez-Quinto F, Santpere G, Chiang CWK, DeGiorgio M,

486        Prado-Martinez J, Rodriguez JA, Rasmussen S, Quilez J, Ramirez O,

487        Marigorta UM et al. 2014. Derived immune and ancestral pigmentation alleles

488        in a 7,000-year-old Mesolithic European. Nature 507(7491):225-228.

489    Perry GH, Dominy NJ, Claw KG, Lee AS, Fiegler H, Redon R, Werner J, Villanea

490        FA, Mountain JL, Misra R, Carter NP, Lee C et al. 2007. Diet and the

491        evolution of human amylase gene copy number variation. Nat Genet

492        39(10):1256-1260.

493    Perry GH, Tchinda J, McGrath SD, Zhang JJ, Picker SR, Caceres AM, Iafrate AJ,

494        Tyler-Smith C, Scherer SW, Eichler EE, Stone AC, Lee C. 2006. Hotspots for

495        copy number variation in chimpanzees and humans. Proc Natl Acad Sci USA

496        103(21):8006-8011.

497    Pronk JC, Frants RR. 1979. New genetic variants of parotid salivary amylase. Hum

498        Hered 29(3):181-186.

499    Pronk JC, Frants RR, Jansen W, Eriksson AW, Tonino GJM. 1982. Evidence for

500        duplication of the human salivary amylase gene. Hum Genet 60(1):32-35.

501    Prufer K, Racimo F, Patterson N, Jay F, Sankararaman S, Sawyer S, Heinze A,

502        Renaud G, Sudmant PH, de Filippo C, Li H, Mallick S et al. 2014. The

503       complete genome sequence of a Neanderthal from the Altai Mountains.

504       Nature 505(7481):43-49.

505  Sudmant PH, Kitzman JO, Antonacci F, Alkan C, Malig M, Tsalenko A, Sampas N,

506       Bruhn L, Shendure J, Eichler EE, Project G. 2010. Diversity of Human Copy

507       Number Variation and Multicopy Genes. Science 330(6004):641-646.

508  Usher CL, Handsaker RE, Esko T, Tuke MA, Weedon MN, Hastie AR, Cao H, Moon

509       JE, Kashin S, Fuchsberger C, Metspalu A, Pato CN et al. 2015. Structural

510       forms of the human amylase locus and their relationships to SNPs,

511       haplotypes and obesity. Nat Genet 47(8):921-925.

512  Yong RYY, Mustaffa SAB, Wasan PS, Sheng L, Marshall CR, Scherer SW, Teo Y-Y,

513       Yap EPH. 2016. Complex Copy Number Variation Of Amy1 Does not

514       Associate With Obesity in Two East Asian Cohorts. Hum Mutat 37(7):669-

515       678.

516  Yuan B, Liu P, Gupta A, Beck CR, Tejomurtula A, Campbell IM, Gambin T, Simmons

517       AD, Withers MA, Harris RA, Rogers J, Schwartz DC et al. 2015. Comparative

518       Genomic Analyses of the Human NPHP1 Locus Reveal Complex Genomic

519       Architecture and Its Regional Evolution in Primates. PLoS Genet

520       11(12):e1005686.

521

522

523

# Figure legends

**Figure 1. Structures of the reference allele and an allele carrying the *AMY2A+2B* duplication**

Most alleles at the human amylase locus conform to the general structure exemplified by the sequence in the human reference assembly (AMY2B_011111, upper diagram), with one copy each of *AMY2A* and *AMY2B*, and an odd number of copies of *AMY1* (in this case 3). Other members of this allelic series, with odd numbers of copies of *AMY1*, have different numbers (including none) of the 95kb unit shown, containing two copies of *AMY1* and the *AMY2A* pseudogene designated *AMYP1*. The lower diagram shows, on the same scale, the simplest example of a structure containing the *AMY2A+2B* duplication (AMY2B_022101); duplication of a 116kb sequence encompassing *AMY2B*, *AMY2A* and one copy of *AMY1* leads to the formation of a haplotype with 2 copies each of *AMY2B*, *AMY2A* and *AMY1*. Other members of this same allelic series can contain higher even numbers of *AMY1*, again differing in numbers of the 95kb (*AMY1*)$_2$–*AMYP1* unit shown above. Note that the non-homologous duplication is accompanied by the formation of a specific sequence junction between sequences upstream of *AMY2B* and the 18kb repeat sequence between *AMY1A* and *AMY1B*, indicated here by "J" [Carpenter et al., 2015; Usher et al., 2015].

**Figure 2. New experimental methods for high-resolution measurement of *AMY1* copy number and *AMY2A/2B* duplication**

24

547 (a) Results from 854 *AMY1* copy number assays, each with two measurements of

548 *AMY1* copy number using the NF2 and NF5 variants of the *AMY1*C PRT assay (see

549 Materials and Methods). The appearance of clear clusters allows the confident

550 assignment of nearly all samples to integer copy numbers based on these two

551 PCRs, especially at copy numbers below 10. (b) Quantification of the junction

552 fragment for the *AMY2A/2B* duplication allele and its derivatives, measuring the

553 representation of a PCR product from the specific duplication junction fragment

554 ("dup") relative to a control product present in two copies in every individual. Traces

555 are shown from NA18854, NA18859, NA19116 and NA19200, with 0, 1, 2 and 3

556 copies of the duplication junction respectively.

557

558 **Figure 3. Segregation of amylase haplotypes in family trio Y060 demonstrated**

559 **by microsatellite and fibre-FISH analysis**

560 The first table summarises measured copy numbers for *AMY1*, *AMY2A*, *AMY2B* and

561 the junction sequence in this family (see also Table 3). The microsatellite allele

562 profiles demonstrate the split of the total *AMY1* copy number between the different

563 allele lengths (for example, the 12 copies of the father NA19119 are split 1 + 4 + 7).

564 There are four possible segregation patterns for this trio logically compatible with the

565 total copy numbers and whole-number splits. The untransmitted allele in the mother

566 NA19116 carries one copy each of *AMY2A* and *AMY2B*, and is therefore strongly

567 predicted to have an odd number of copies of *AMY1*. Only two of the four possible

568 segregation patterns have an odd number of *AMY1* in the untransmitted maternal

569 allele, and both of those involve transmission of 3 copies from the mother NA19116,

570 and 8 copies from the father NA19119 (Table 3). One of those compatible

571  segregation patterns is indicated here by the arrows and numbers. These analyses

572  together suggest the haplotype segregation shown in the lower table, with

573  transmitted alleles shown in orange (AMY2B_012341, paternal) and blue

574  (AMY2B_033201, maternal).

575  In fibre-FISH analysis, the *AMY2B* probe employed (green) is specific to sequence

576  upstream of *AMY2B*. The probe (red) for the sequence upstream of *AMY2A* cross-

577  hybridizes with very similar sequence surrounding the ERV upstream of *AMY1,* and

578  the *AMY1* gene probe (white) also cross-hybridizes with coding regions of *AMY2A*

579  and *AMY2B*. In many locations, this additional cross-hybridization between similar

580  amylase sequences provides useful confirmation of the type and orientation of the

581  gene. Examples of hybridization observed with these three probes with *AMY1*,

582  *AMY2A* and *AMY2B* are shown in the top panel. The orange box frames images

583  from the haplotype AMY2B_012341 transmitted from father to child, with the

584  composition $(AMY1)_8(AMY2A)_2(AMY2B)_1$, including a duplicated copy of *AMY2A* in

585  the forward orientation preceded by 3 copies of *AMY1*. The reconstructed

586  interpretation of the ≈490kb structure appears to be seamless, in that it includes no

587  new short-range junctions, but the overall arrangement suggests that it arose

588  independently of the untransmitted paternal $(AMY1)_4(AMY2A)_2(AMY2B)_1$ allele

589  (AMY2B_012211) shown in Supp. Figure S9. In the blue box, full-length haplotype

590  images from the triplication allele (AMY2B_033201) transmitted from mother to child

591  (Table 3) are shown above the inferred gene arrangement, and finally the full

592  (≈300kb) haplotype reconstruction. "J" shows the inferred positions of the duplication

593  junction sequence, and the boxed region highlights the inversion of one copy each of

594  *AMY1* and *AMY2A* relative to the reference assembly orientation.
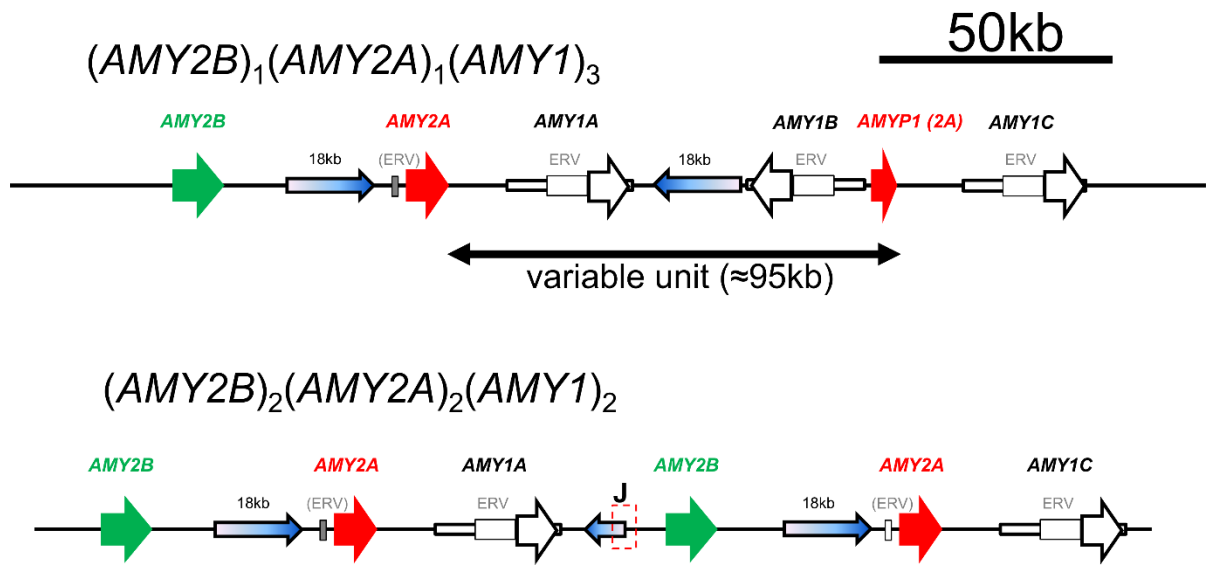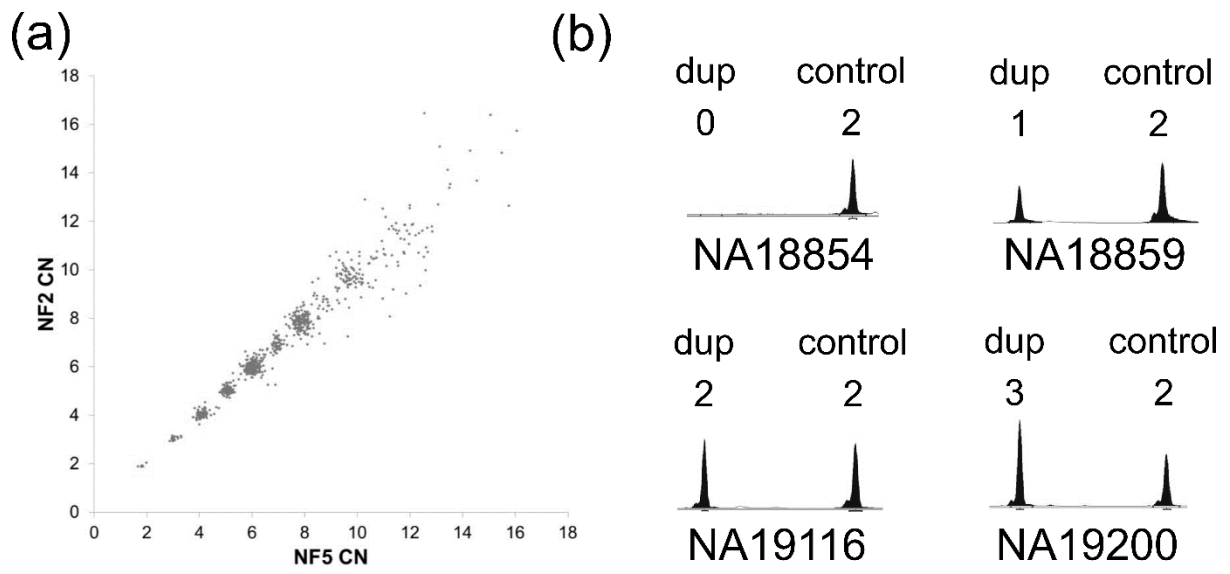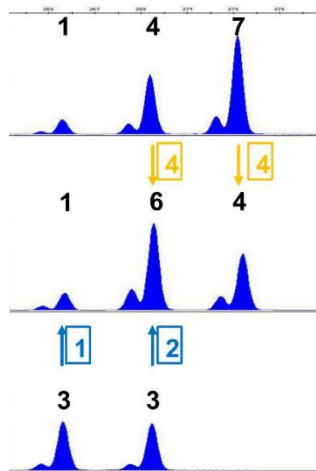
595

Figure 1



$(AMY2B)_1(AMY2A)_1(AMY1)_3$

$(AMY2B)_2(AMY2A)_2(AMY1)_2$

Figure 2

(a)



(b)

Figure 3

Diploid integer copy numbers:

| component | AMY1 | AMY2A | AMY2B | junction |
|---|---|---|---|---|
| Father (NA19119) | 12 | 4 | 2 | 0 |
| Child (NA19120) | 11 | 5 | 4 | 2 |
| Mother (NA19116) | 6 | 4 | 4 | 2 |

**Hybridization patterns**



NA19119 (father)
AMY1 total: 12

8 copies

NA19120 (child)
AMY1 total: 11

3 copies

NA19116 (mother)
AMY1 total: 6

Use trio segregation with microsatellite profiles to deduce haplotype copy numbers

Haplotype copy numbers:

| ID | component | AMY1 | AMY2A | AMY2B | junction |
|---|---|---|---|---|---|
| Father | transmitted | 8 | 2 | 1 | 0 |
| | untransmitted | 4 | 2 | 1 | 0 |
| Child | Paternal | 8 | 2 | 1 | 0 |
| | Maternal | 3 | 3 | 3 | 2 |
| Mother | transmitted | 3 | 3 | 3 | 2 |
| | untransmitted | 3 | 1 | 1 | 0 |

Father GM19119: $(AMY1)_8(AMY2A)_2(AMY2B)_1$ ID:AMY2B_012341

Child GM19120: $(AMY1)_8(AMY2A)_2(AMY2B)_1$ ID:AMY2B_012341

Mother GM19116: $(AMY1)_3(AMY2A)_3(AMY2B)_3$ ID:AMY2B_033201

Child GM19120: $(AMY1)_3(AMY2A)_3(AMY2B)_3$ ID:AMY2B_033201

29

**Table 1. Primers used in this work.**

| Primer | Sequence (5'-3') |
|---|---|
| AMY1CF | TTCTAAGGTGCCTTCTAGTC |
| AMY1CR | CATCTTCAAGCCTGCATTC |
| NF2 | ATAGCTTAGAGTAGTTAAC |
| AMY1CRB2 | AGTGAGATGAGGCATTGTG |
| NF5 | GGCCTCTATACATGAG |
| AMY2B2D | GCCTGGCTAATTTGTTGTTAG |
| AMY2B2R | AAATTAACTCCATGCATCACC |
| AMY2B2F | TGCATAGAAATGGCACATAGT |
| AMY1_2F | ACAGTTGATTTTTGATCTTGTAGG |
| AMY1_2R | TACAGCATCCACATAAATACGAA |

**Table 2. Segregation of *AMY1*, *AMY2A* and *AMY2B* copy number in Yoruban trios Y045 and Y056**

| Family Y045 | ID | component | *AMY1* | *AMY2A* | *AMY2B* | junction |
|---|---|---|---|---|---|---|
| Mother | NA19201 | diploid | 6 | 2 | 2 | 0 |
| | | transmitted haplotype | 3 | 1 | 1 | 0 |
| | | untransmitted haplotype | **3** | **1** | **1** | **0** |
| Father | NA19200 | diploid | 7 | 5 | 5 | 3 |
| | | transmitted haplotype | **4** | **4** | **4** | **3** |
| | | untransmitted haplotype | 3 | 1 | 1 | 0 |
| Child | NA19202 | diploid | 7 | 5 | 5 | 3 |
| | | Maternal haplotype | **3** | **1** | **1** | **0** |
| | | Paternal haplotype | **4** | **4** | **4** | **3** |
| | | | | | | |
| Family Y056 | ID | component | *AMY1* | *AMY2A* | *AMY2B* | junction |
| Mother | NA19159 | diploid | 8 | 6 | 6 | 4 |
| | | transmitted haplotype | **3** | **1** | **1** | **0** |
| | | untransmitted haplotype | 5 | 5 | 5 | 4 |
| Father | NA19160 | diploid | 6 | 2 | 2 | 0 |
| | | transmitted haplotype | **3** | **1** | **1** | **0** |
| | | untransmitted haplotype | 3 | 1 | 1 | 0 |
| Child | NA19161 | diploid | 6 | 2 | 2 | 0 |
| | | Maternal haplotype | **3** | **1** | **1** | **0** |
| | | Paternal haplotype | **3** | **1** | **1** | **0** |

**Table 3. Segregation of *AMY1*, *AMY2A* and *AMY2B* copy number in Yoruban trios Y060 and Y072**

| Family Y060 | ID | component | *AMY1* | *AMY2A* | *AMY2B* | junction |
|---|---|---|---|---|---|---|
| Mother | NA19116 | diploid | 6 | 4 | 4 | 2 |
| | | transmitted haplotype | **3** | **3** | **3** | **2** |
| | | untransmitted haplotype | 3 | 1 | 1 | 0 |
| Father | NA19119 | diploid | 12 | 4 | 2 | 0 |
| | | transmitted haplotype | **8** | **2** | **1** | **0** |
| | | untransmitted haplotype | 4 | 2 | 1 | 0 |
| Child | NA19120 | diploid | 11 | 5 | 4 | 2 |
| | | Maternal haplotype | **3** | **3** | **3** | **2** |
| | | Paternal haplotype | **8** | **2** | **1** | **0** |
| | | | | | | |
| Family Y072 | ID | component | *AMY1ᵃ* | *AMY2A* | *AMY2B* | junction |
| Mother | NA19152 | diploid | 5/6 | 6 | 5 | 3 |
| | | transmitted haplotype | **2/3** | **3** | **2** | **1** |
| | | untransmitted haplotype | 3 | 3 | 3 | 2 |
| Father | NA19153 | diploid | 8 | 2 | 2 | 0 |
| | | transmitted haplotype | **3** | **1** | **1** | **0** |
| | | untransmitted haplotype | 5 | 1 | 1 | 0 |
| Child | NA19154 | diploid | 5/6 | 4 | 3 | 1 |
| | | Maternal haplotype | **2/3** | **3** | **2** | **1** |
| | | Paternal haplotype | **3** | **1** | **1** | **0** |

a Alternative values are shown for the *AMY1* copy numbers of the mother and child in family Y072; because of the partial copy of *AMY1* on the transmitted maternal haplotype, the copy number recorded depends on the precise location of the measure used.