

CONGA: Copy number variation genotyping in ancient genomes and low-coverage sequencing data

Arda Söylev^{1*}, Sevim Seda Çokoglu², Dilek Koptekin³, Can Alkan⁴, and Mehmet Somel^{2*}

¹*Department of Computer Engineering, Konya Food and Agriculture University, Konya, 42080, Turkey*

²*Department of Biology, Middle East Technical University, Ankara, 06800, Turkey*

³*Department of Health Informatics, Middle East Technical University, Ankara, 06800, Turkey*

⁴*Department of Computer Engineering, Bilkent University, Ankara, 06800, Turkey*

*Corresponding authors: arda.soylev@gidatarim.edu.tr and msomel@metu.edu.tr

ABSTRACT

1 To date, ancient genome analyses have been largely confined to the study of single nucleotide
2 polymorphisms (SNPs). Copy number variants (CNVs) are a major contributor of disease and of
3 evolutionary adaptation, but identifying CNVs in ancient shotgun-sequenced genomes is hampered
4 by (a) most published genomes being $<1\times$ coverage, (ii) ancient DNA fragments being typically
5 <80 bps. These characteristics preclude state-of-the-art CNV detection software to be effectively
6 applied to ancient genomes. Here we present CONGA, an algorithm tailored for genotyping deletion
7 and duplication events in genomes with low depths of coverage. Simulations show that CONGA can
8 genotype deletions and duplications >1 Kbps with F-scores >0.77 and >0.82 at $\geq 0.5\times$, respectively.
9 Further, down-sampling experiments using published ancient BAM files reveal that >1 Kbps deletions
10 could be genotyped at F-score >0.75 at $\geq 1\times$ coverage. Using CONGA, we analyse deletion events at
11 10,018 loci in 56 ancient human genomes spanning the last 50,000 years, with coverages $0.4\times$ - $26\times$.
12 We find inter-individual genetic diversity measured using deletions and SNPs to be highly correlated,
13 suggesting that deletion frequencies broadly reflect demographic history. We also identify signatures
14 of purifying selection on deletions, such as an excess of singletons compared to those in SNPs.
15 CONGA paves the way for systematic studies of drift, mutation load, and adaptation in ancient and
16 modern-day gene pools through the lens of CNVs.

17 **Keywords** Genomics · ancient DNA · CNV genotyping · deletion · low coverage whole genome sequencing

18 1 Introduction

19 Ancient genomics, the analysis of genetic material extracted from archaeological and paleontological remains, has
20 become a major source of information for the study of population history and evolution over the last decade (Skoglund
21 and Mathieson, 2018; Frantz *et al.*, 2020; Shapiro and Hofreiter, 2014; Marciniak and Perry, 2017). While the number of
22 published ancient genomes is exponentially growing, their analyses have yet been nearly exclusively limited to those of
23 single-nucleotide polymorphisms (SNPs), while structural variations (SVs) in ancient genomes remain mostly ignored.
24 Copy number variations (CNVs) are a common type of SVs and include deletions and duplications ranging from 50
25 bps to several megabasepairs. Although their number, by count, is much fewer than SNPs, the fraction of the genome
26 affected by CNVs is well past that accounted for SNPs (Conrad *et al.*, 2010). Likewise, CNVs are a major contributor
27 to phenotypic variation: they are frequently discovered as the basis of diverse biological adaptations (Gonzalez *et al.*,
28 2005; Perry *et al.*, 2007; Xue *et al.*, 2008; Chan *et al.*, 2010; McLean *et al.*, 2011; Hardwick *et al.*, 2011; Kothapalli
29 *et al.*, 2016; Nuttle *et al.*, 2016; Hsieh *et al.*, 2019) as well as genetic diseases (reviewed in (Zhang *et al.*, 2009; Saitou
30 and Gokcumen, 2020; Stankiewicz and Lupski, 2010; Girirajan *et al.*, 2011)). This renders the study of CNVs in
31 ancient genomes two-fold attractive. First, as CNVs frequently serve as genetic material for adaptation, their study
32 in ancient genomes can allow detailed temporal investigation of adaptive processes. Examples include evolutionary
33 changes in salivary amylase copy numbers in humans and in dogs, thought to represent responses to a shift to starch-rich
34 diets (Mathieson and Mathieson, 2018; Bergström *et al.*, 2020). Second, large deletions can be a major source of
35 deleterious mutation load, and studying deletion frequencies in ancient genome samples from extinct species or severely
36 bottlenecked populations can be highly revealing about the genetic health of lineages. For instance, a study on the last

37 surviving mammoth population on Wrangel Island has found an excess of deletions in this sample, which could have
38 compromised the population's fitness (Rogers and Slatkin, 2017).

39 Despite this appeal, the impact of CNVs on evolutionary history and ancient phenotypes remains largely unex-
40 plored (Frantz *et al.*, 2020). The reason lies in the significant technical challenges in CNV detection posed by ancient
41 genomes. State-of-the-art methods for CNV discovery from shotgun genome sequencing data require at least moderate
42 depth of coverage (Abyzov *et al.*, 2011; Boeva *et al.*, 2012; Alkan, 2020; Smith *et al.*, 2015) and read-pair informa-
43 tion (Soylev *et al.*, 2017, 2019; Rausch *et al.*, 2012; Layer *et al.*, 2014; Chen *et al.*, 2016; Einfeldt *et al.*, 2017), or long
44 reads (Sedlazeck *et al.*, 2018; Chaisson *et al.*, 2015). However, due to the degraded and elusive nature of ancient DNA,
45 ancient genome data is frequently produced at low coverage ($<1\times$) and the molecules retrieved are typically short,
46 between 50-80 bps. Available CNV discovery methods are therefore inapplicable to most ancient genome data sets, and
47 so far, no specific algorithm for CNV identification in ancient genomes has been developed.

48 With the aim to fill this gap, here we present CONGA (Copy Number Variation Genotyping in Ancient Genomes and
49 Low-coverage Sequencing Data), a CNV genotyping algorithm tailored for ancient genomes, which estimates copy
50 number beyond presence/absence of events. To our knowledge, CONGA is the first tool specifically developed for this
51 purpose. There are four reasons behind the use of CNV genotyping instead of CNV discovery: (a) CNV discovery using
52 low coverage ancient genomes is impractical if not impossible; (b) for many species studied using ancient genomics,
53 CNV reference sets based on high quality genomes are already available; (c) ancient variation will largely overlap with
54 modern-day variation in the vast majority of cases; (d) genotyping has much shorter running times and lower memory
55 usage than discovery. Indeed, most human ancient genome studies to date have chosen genotyping over *de novo* SNP
56 discovery (Prüfer, 2018; Link *et al.*, 2017). It should likewise be possible to genotype CNVs in low coverage genomes
57 with high accuracy and in short running times using depth of coverage and split-read information. We therefore believe
58 that CONGA's approach can efficiently open up CNV analyses to ancient genome studies, as well as low coverage
59 whole genome sequencing datasets of extant organisms.

60 We first test the accuracy of CONGA using purely simulated ancient genomes as well as down-sampling simulations of
61 real ancient genomes. We show that our algorithm exhibits reliable performance for deletions and duplications even at
62 low depths of coverage (i.e., $<1\times$). We also discover that published ancient genome BAM files may frequently not
63 be suitable for duplication genotyping at low coverage. We next genotype deletions in a set of real ancient human
64 genomes using deletions ascertained in present-day African genomes (Figure 1). We find that accurate CNV genotyping
65 is sensitive to a number of technical parameters, requiring comprehensive filtering to achieve a high quality CNV call
66 set. Having created such a deletion CNV set, we show that CNV diversity largely represents demographic history,
67 paralleling SNP diversity among the same ancient individuals. We further present evidence for negative selection on
68 deletions in these ancient genomes.

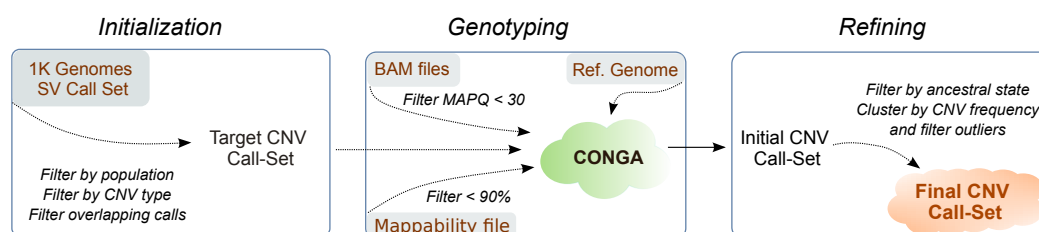


Figure 1: Overall workflow of CONGA. The first step involves initialization, where we create the input (reference) CNV file using the deletions and duplications of African populations (AFR) from Phase 3 of the 1000 Genomes Project SV call set (Sudmant *et al.*, 2015). We apply our genotyping algorithm in the second step and create the initial CNV call set. We then perform a filtering and refining step, which is used to generate the final CNV call set.

69 2 Results

70 2.1 Overview of the algorithm

71 We developed CONGA to genotype given candidate CNVs in mapped read (BAM) files (Methods). Briefly, CONGA
72 first calculates the number of reads mapped to each given interval in the reference genome, which we call “observed
73 read-depth”. It then calculates the “expected diploid read-depth”, i.e., the GC-content normalized read-depth given the
74 genome average. Using these values, CONGA calculates the likelihood for each genotype by modeling the read-depth

75 distribution as Poisson, similar to common CNV callers. The genotypes can be homozygous CNV, heterozygous
76 CNV, or no CNV. Using these likelihoods CONGA then calculates a statistic that we term the C-score, defined as the
77 likelihood of a CNV being true (in heterozygous or homozygous state) over it being false (no CNV). For genotyping
78 duplications, CONGA also uses an additional split-read step in order to utilize paired-end information. Briefly, it splits
79 reads and remaps the split within the genome, treating the two segments as paired-end reads (Karakoc *et al.*, 2012;
80 Soylev *et al.*, 2019). Either type of signature, read-depth or paired-end, can be sufficient to call a duplication (Methods).

81 **2.2 Accuracy evaluation using simulated genomes and comparison with CNV discovery tools**

82 To evaluate the performance of our algorithm we first simulated genomes with CNVs of ancient-like characteristics. We
83 employed VarSim (Mu *et al.*, 2015) to insert deletions and duplications into the human reference genome (GRCh37).
84 We used three different size intervals for CNVs: small (100 bps - 1000 bps), medium (1,000 bps - 10,000 bps) and
85 large (10,000 bps - 100,000 bps). We thus simulated three genomes, each with roughly 1,500 deletions and 1,500
86 duplications of a specific size range (see Supplemental Fig. 1 for the exact numbers and length distributions of CNVs
87 inserted in each genome). We next used these genomes as input to the ancient read simulator Gargammel (Renaud
88 *et al.*, 2017), which generates paired-end short Illumina reads with varying fragment sizes (median 66 bps) as well as
89 post-mortem damage. The data was generated at various depths: 0.05 \times , 0.1 \times , 0.5 \times , 1 \times and 5 \times (Methods). We then
90 used CONGA to genotype CNVs across the simulated ancient genomes using the candidate CNV call set. In order
91 to assess specificity in addition to sensitivity, we also used a background (false) CNV list, prepared using published
92 deletion and duplication calls from modern-day human long read sequencing datasets (Audano *et al.*, 2019; Chaisson
93 *et al.*, 2019; Zook *et al.*, 2020; Collins *et al.*, 2020), as well as from African populations (AFR) from Phase 3 of the
94 1000 Genomes Project (Sudmant *et al.*, 2015). We mixed these false CNVs to the list of true CNVs with a ratio of
95 approximately 10:1 (roughly 15,000 false events vs. 1,500 true events), and used this mixed list as the candidate CNV
96 call set to CONGA (Methods). To assess the performance of CONGA in identifying CNVs, we further compared it with
97 two of the widely used CNV discovery approaches: CNVnator (Abyzov *et al.*, 2011) and FREEC (Boeva *et al.*, 2012).

98 Table 1 shows true and false predictions by CONGA, FREEC and CNVnator, as well as their true positive rate (TPR),
99 false discovery rate (FDR) and the F-Score for identifying deletions and duplications of small, medium and large size
100 (as defined above). In detecting deletions and duplications, CONGA achieved a substantially lower FDR and overall
101 higher sensitivity than the other two methods. Note that even at 0.05 \times coverage, CONGA yields accurate results for
102 medium and large sized CNVs, whereas the other tools clearly fail. On the other hand, the performance of each tool
103 becomes comparable as the coverage approaches to depths of 5 \times with large CNV sizes. For small CNVs (<1 Kbps), all
104 three tools underperform, although CONGA predictions for small deletions still have higher recall and precision than
105 the other two tools (see Supplemental Fig. 2 for precision-recall curves). We note that CONGA's genotyping approach
106 naturally boosts its performance over FREEC and CNVnator, which are CNV discovery tools. Still, the fact that our
107 candidate CNV call set included 10 times more false CNVs than true CNVs indicates the overall reliability of CONGA
108 in CNV identification in ancient genomes relative to these alternatives.

109 The simulation results thus suggest that CONGA can efficiently and accurately genotype deletions and duplications of
110 length >1 Kbps in ancient genomes at $\geq 0.5\times$ coverage. We then turned to studying the performance of CONGA in
111 identifying copy numbers for medium and large size CNVs.

112 **2.2.1 Copy number prediction of CNVs**

113 Beyond the identification of deletion and duplication events, classifying individual genotypes as heterozygous or
114 homozygous CNVs could provide valuable information for population genetic analyses of CNVs. However, predicting
115 CNV copy numbers can be a significant challenge on low coverage ancient genomes (Kousathanas *et al.*, 2017). We
116 thus assessed the performance of CONGA to determine the copy number of a CNV based on the likelihood model
117 described above using our simulation data. For simplicity, CONGA only evaluates the possibility of homozygous
118 duplications (ignoring copy numbers ≥ 3). Figure 2 shows CONGA's copy number prediction performance for deletions
119 and duplications using F-scores for each coverage tested. We found that F-scores were above ≥ 0.5 at coverages $\geq 0.5\times$.
120 Importantly, CONGA had comparable power in identifying heterozygous and homozygous events of size >1 Kbps
121 (Supplemental Table S1.B).

122 **2.3 Down-sampling experiments with real ancient genomes**

123 Beyond pure simulations, we further studied the performance of CONGA in identifying CNVs at various depths
124 of coverage using real ancient genome data. As no ground truth CNV call-set is available, we used the following
125 approach: (i) we chose three published ancient genomes (BAM or FASTQ files) of relatively high coverage ($\geq 9\times$), (ii)
126 we genotyped CNVs using the full genome data with CONGA and with a modern-day human CNV call set, (iii) we

Table 1: Summary of simulation predictions by CONGA, FREEC and CNVnator.

	Cov.	Total	CONGA					FREEC					CNVnator				
			True	False	TPR	FDR	F-Score	True	False	TPR	FDR	F-Score	True	False	TPR	FDR	F-Score
Dels (small)	0.05×	1810	1474	3724	0.81	0.72	0.42	0	1221	0.00	1.00	-	3	47442	0.00	1.00	0.00
	0.1×	1810	1268	2415	0.70	0.66	0.46	0	198	0.00	1.00	-	0	402	0.00	1.00	-
	0.5×	1810	767	456	0.42	0.37	0.51	0	6761	0.00	1.00	-	0	806	0.00	1.00	-
	1×	1810	834	189	0.46	0.18	0.59	0	1916	0.00	1.00	-	0	263	0.00	1.00	-
	5×	1810	1236	89	0.68	0.07	0.79	20	392	0.01	0.95	0.02	341	493	0.19	0.59	0.26
Dups (small)	0.05×	1751	230	1878	0.13	0.89	0.12	0	44	0.00	1.00	-	7	47700	0.00	1.00	0.00
	0.1×	1751	128	1116	0.07	0.90	0.09	0	7	0.00	1.00	-	0	28699	0.00	1.00	-
	0.5×	1751	304	28	0.17	0.08	0.29	0	3	0.00	1.00	-	0	9	0.00	1.00	-
	1×	1751	618	22	0.35	0.03	0.52	0	555	0.00	1.00	-	0	884	0.00	1.00	-
	5×	1751	1183	30	0.68	0.02	0.80	35	77	0.02	0.69	0.04	2	0	0.00	0.00	0.00
Dels (medium)	0.05×	1680	868	791	0.52	0.48	0.52	0	83	0.00	1.00	-	0	68	0.00	1.00	-
	0.1×	1680	821	328	0.49	0.29	0.58	0	237	0.00	1.00	-	1	216	0.00	1.00	0.00
	0.5×	1680	1066	68	0.63	0.06	0.76	239	6433	0.14	0.96	0.06	187	257	0.11	0.58	0.18
	1×	1680	1268	56	0.75	0.04	0.84	421	2135	0.25	0.84	0.20	330	257	0.20	0.44	0.29
	5×	1680	1380	53	0.82	0.04	0.89	929	485	0.55	0.34	0.60	949	423	0.56	0.31	0.62
Dups (medium)	0.05×	1684	193	29	0.11	0.13	0.20	0	3	0.00	1.00	-	0	114	0.00	1.00	-
	0.1×	1684	384	20	0.23	0.05	0.37	0	3	0.00	1.00	-	0	102	0.00	1.00	-
	0.5×	1684	990	10	0.59	0.01	0.74	271	15	0.16	0.05	0.28	2	4	0.00	0.67	0.00
	1×	1684	1233	15	0.73	0.01	0.84	582	937	0.35	0.62	0.36	16	2	0.01	0.11	0.02
	5×	1684	1507	13	0.89	0.01	0.94	1000	329	0.59	0.25	0.66	105	2	0.06	0.02	0.12
Dels (large)	0.05×	1385	931	1017	0.67	0.52	0.56	0	87	0.00	1.00	-	84	131	0.06	0.61	0.11
	0.1×	1385	1025	549	0.74	0.35	0.69	0	754	0.00	1.00	-	560	246	0.40	0.31	0.51
	0.5×	1385	1093	334	0.79	0.23	0.78	664	3136	0.48	0.83	0.26	1049	293	0.76	0.22	0.77
	1×	1385	1083	373	0.78	0.26	0.76	1239	156	0.89	0.11	0.89	1204	309	0.87	0.20	0.83
	5×	1385	1072	383	0.77	0.26	0.75	1260	154	0.91	0.11	0.90	1265	453	0.91	0.26	0.82
Dups (large)	0.05×	1532	966	116	0.63	0.11	0.74	0	6	0.00	1.00	-	4	354	0.00	0.99	0.00
	0.1×	1532	1157	138	0.76	0.11	0.82	0	0	-	-	-	455	315	0.30	0.41	0.40
	0.5×	1532	1388	205	0.91	0.13	0.89	589	97	0.38	0.14	0.53	1039	77	0.68	0.07	0.78
	1×	1532	1402	223	0.92	0.14	0.89	1305	266	0.85	0.17	0.84	1216	94	0.79	0.07	0.86
	5×	1532	1410	243	0.92	0.15	0.89	1304	294	0.85	0.18	0.83	1350	165	0.88	0.11	0.89

Table shows the CNV prediction performance of CONGA, FREEC and CNVnator on simulated genomes with depth 0.05×, 0.1×, 0.5×, 1× and 5× for deletions (Dels) and duplications (Dups) of multiple CNV size intervals including 100 bps - 1 Kbps (small), 1 Kbps - 10 Kbps (medium) and 10 Kbps - 100 Kbps (large). Here, **True** and **False** refer to correct and incorrect predictions respectively, **TPR** is true positive rate (or recall) and **FDR** is false discovery rate ($1 - Precision$) of each algorithm. The **F-Score**, is calculated as $(2 \times Precision \times Recall) / (Precision + Recall)$. Bold values in each row represent the highest TPR, lowest FDR, or highest F-Score across the tools. CONGA consistently shows the best performance in identifying small and medium sized deletions and duplications, as well as the best performance in identifying large CNVs at low coverage. See Supplemental Material for the commands used to run each tool and Supplemental Table S1.A for details.

127 down-sampled the ancient genome data to lower coverages, (iv) we assessed CONGA's performance in genotyping the
128 same CNVs at low coverage (Methods).

129 Specifically, we selected a ($\sim 23.3\times$) ancient Eurasian genome (Yamnaya) (de Barros Damgaard *et al.*, 2018b), a
130 $13.1\times$ ancient genome from Greenland (Saqqaq) (Rasmussen *et al.*, 2010), and a $9.6\times$ ancient genome from Ethiopia
131 (Mota) (Llorente *et al.*, 2015). The Yamnaya genome was only available as a BAM file, while the latter two were
132 available as FASTQ files, which we processed into BAM files (Methods). We used a list of modern-day human CNVs
133 as candidate CNV set ($n = 17,392$ deletions and $n = 14,888$ duplications) (Methods) as input to CONGA. We thus
134 genotyped between 688-1581 deletions and 638-4097 duplications across these three genomes using the full data. We
135 then down-sampled all three BAM files to various depths, and repeated the genotyping for each genome. We estimated
136 the CONGA's TPR (true positive rate) and FDR (false discovery rate) on down-sampled genomes by treating the CNVs
137 genotyped using the full data as ground truth (Methods).

138 CONGA displayed satisfactory performance in identifying deletions in all three genomes, with TPR of $>70\%$ and
139 FDR of $<45\%$ at coverages around $0.5\times$ (Figure 3, Supplemental Table S1.C). For duplications, however, CONGA's

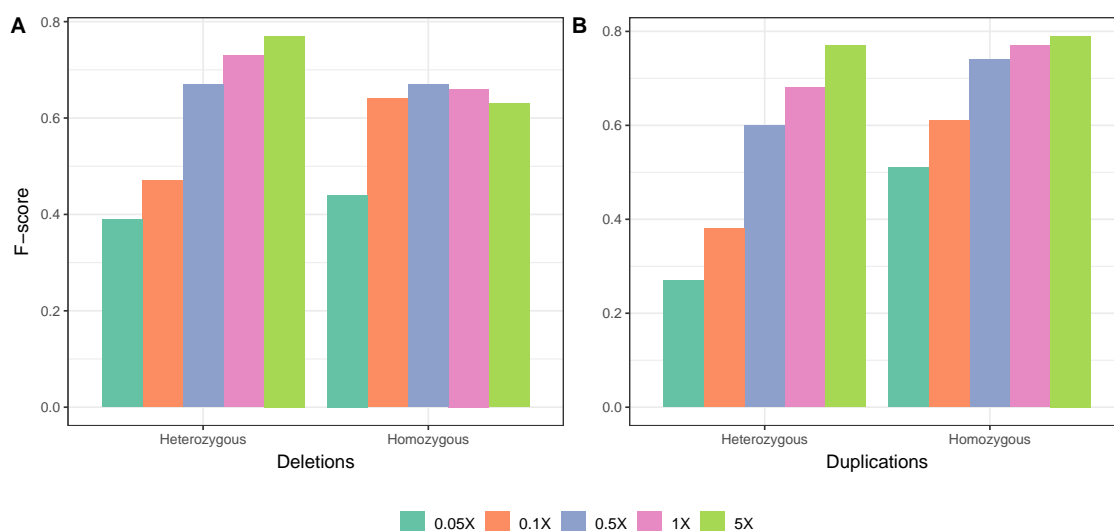


Figure 2: Copy-number prediction performance of CONGA using F-score for (A) deletions and (B) duplications using medium and large CNVs merged for various coverage values.

140 performance was lower. At around $1\times$ coverage, duplication TPR was $>40\%$ in the Saqqaq and Mota genomes, and
 141 only 22% in the Yamnaya genome.

142 To understand the reasons behind CONGA's striking under-performance at low coverages, especially on the Yamnaya
 143 genome, we studied the results in further detail. We noticed that only 47 (2.8%) of the 1,661 originally called
 144 duplication events in the full Yamnaya genome were genotyped using read-depth information (C-score), while the rest
 145 were genotyped using paired-end read information. This can explain CONGA's under-performance at low coverage in
 146 the Yamnaya genome: as coverage decreases, the number of paired-end reads supporting a duplication falls rapidly,
 147 compromising recall. The lack of usable read-depth information in the full Yamnaya genome, in turn, could be standard
 148 quality filters applied to the BAM file before data publication. Such filtering could have erased the read-depth signature,
 149 leaving only paired-end information available, which is not helpful at low coverages. This scenario is supported by the
 150 fact CONGA duplication calls are clearly more successful in the Saqqaq and Mota genomes at low coverage. These two
 151 were retrieved as original FASTQ files instead of processed BAM files. Moreover, 35% (1498/4027) and 27% (172/638)
 152 of duplication events were genotyped using read-depth information on the full Saqqaq and Mota genomes, respectively,
 153 in contrast to only 2.8% in the Yamnaya genome (see Discussion).

154 Overall, both our simulations and down-sampling experiments with real genomes suggest that CONGA can efficiently
 155 genotype >1 Kbps deletion events at depths of coverage of $0.5\times$, and even at $0.1\times$. CONGA could thus be applied
 156 on a large fraction of ancient shotgun sequenced genomes available. In contrast, although CONGA's accuracy in
 157 duplication identification was comparable to those of deletions in simulations, it was dramatically lower in correctly
 158 finding duplications in the down-sampled Yamnaya genome even at $1\times$ coverage. We suspect this low performance is
 159 caused by pre-publication quality filtering of BAM files. Consequently, duplication genotyping with CONGA directly
 160 on published ancient genomes appears infeasible, as ancient genomes are mainly submitted in BAM format in public
 161 repositories (see Discussion). We therefore limited the following analyses on real ancient genomes to deletions >1
 162 Kbps.

163 2.4 Deletion genotyping across ancient genomes reveals strong influence by technical variables

164 We genotyped deletions with CONGA across a diverse sample of real ancient human genomes. Our naive hypothesis
 165 was that CNVs, like SNPs, should display genome-wide similarity patterns that reflect population origin, i.e., shared
 166 genetic drift, among individuals (Conrad and Hurler, 2007; Levy-Sakin *et al.*, 2019; Almarri *et al.*, 2020). We collected
 167 BAM files for 71 ancient human genomes belonging to a time range between c.2,800-45,000 years Before Present (BP)
 168 (Supplemental Table S2) (Rasmussen *et al.*, 2014; Günther *et al.*, 2015; Hofmanová *et al.*, 2016; Jones *et al.*, 2015; Kılınc
 169 *et al.*, 2016; de Barros Damgaard *et al.*, 2018b; Gamba *et al.*, 2014; González-Fortes *et al.*, 2017; de Barros Damgaard
 170 *et al.*, 2018a; Keller *et al.*, 2012; Sikora *et al.*, 2019; Olalde *et al.*, 2014; Lazaridis *et al.*, 2014; Antonio *et al.*, 2019;
 171 Allentoft *et al.*, 2015; Haber *et al.*, 2019; Fu *et al.*, 2014; Broushaki *et al.*, 2016; Seguin-Orlando *et al.*, 2014; Jones

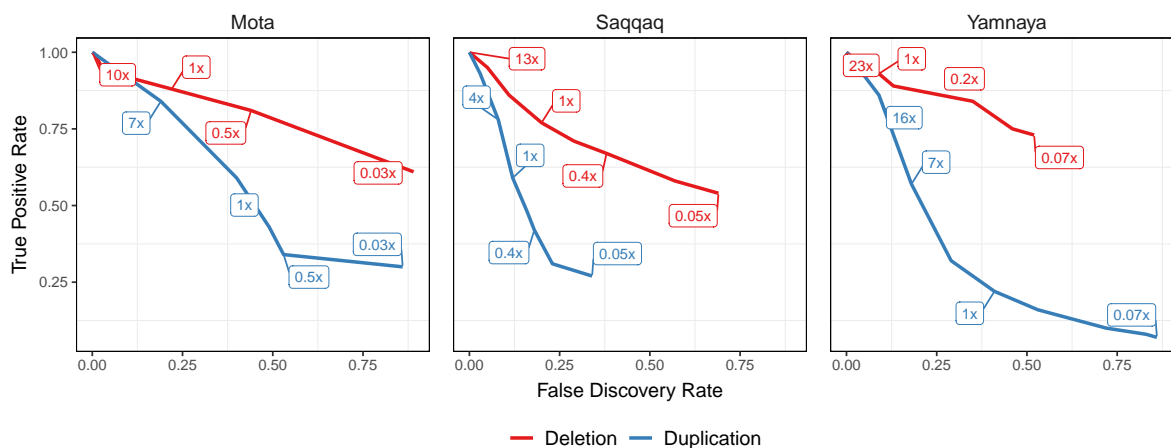


Figure 3: TPR vs FDR curves for deletion and duplication predictions of CONGA using Mota, Saqqaq and Yamnaya genomes down-sampled to various depths from their original coverages of $9.6\times$, $13.1\times$ and $23.3\times$, respectively. The numbers inside boxes show the down-sampled coverage values. We calculated TPR and FDR for down-sampled genomes assuming that our CONGA-based predictions with the original genomes (full data) reflect the ground truth. These predictions, in turn, were made using modern-day CNVs as candidate CNV list. The purpose of the experiment was to evaluate accuracy at lower coverage relative to the full data (Methods).

172 *et al.*, 2017; Haber *et al.*, 2017; Raghavan *et al.*, 2014; Martiniano *et al.*, 2017; Krzewińska *et al.*, 2018; Yaka *et al.*,
 173 2021). These were chosen to bear diverse characteristics, including a wide range in mean coverage ($0.04\times$ - $26\times$,
 174 median = $3.45\times$), population origin (West and East Eurasia and North America), the laboratory of origin (10 different
 175 laboratories), the use of shotgun vs. whole-genome capture protocols, or the use of uracil-DNA-glycosylase (UDG)
 176 treatment (Rohland *et al.*, 2015). For genotyping, we used a candidate CNV dataset of 11,390 autosomal deletions
 177 (>1 Kbps with mean 10,735 bps) identified among African populations (AFR) from Phase 3 of the 1000 Genomes
 178 Project (Sudmant *et al.*, 2015) (Methods). Our motivation for using an African sample here was to avoid ascertainment
 179 bias (Clark *et al.*, 2005) in studying deletion frequencies, as all of the 71 ancient individuals were non-African, and
 180 thus African populations represent an outgroup to our sample set. Genotyping these 11,390 deletions across the 71
 181 BAM files resulted in a median number of 905 deletions [380-3,649] detected per genome, in either heterozygous or
 182 homozygous state. Of the 11,390 deletions, 8,106 (71%) were detected in at least one genome.

183 We then studied deletion copy numbers, or frequencies, across the 71 ancient genomes. Unexpectedly, both hierarchical
 184 clustering (displayed as a heatmap in Supplemental Fig. 3A) and principal components analysis (PCA) (Figure 4A)
 185 revealed that genomes grouped into two in terms of deletion frequencies (Supplemental Fig. 3A). We noticed that
 186 this grouping mainly represented laboratories-of-origin, irrespective of the population-of-origin. One set, which we
 187 call Group 1, contained 30 individuals from 8 laboratories, while the other, Group 2, included 41 genomes from 4
 188 laboratories (Figure 4A). Only 2 of the 10 laboratories had data assigned to both groups. The laboratory-of-origin was
 189 a highly significant factor explaining deletion frequencies in PC1 (Kruskal-Wallis test $p < 10^{-8}$). Further, 735 (9%)
 190 of 8106 deletions identified in at least one genome showed differences in frequency among laboratories (Benjamini-
 191 Hochberg corrected Kruskal-Wallis test $p < 0.05$). Also notably, Group 2 genomes had substantially higher frequencies
 192 of heterozygous deletions: among 735 deletions identified as significantly different among the 71 ancient genomes, the
 193 mean within Group 1 was ~ 4 , while that for Group 2 was ~ 28 .

194 We considered it unlikely that the observed dual grouping could have a biological reason and sought for technical
 195 explanations. To investigate this, we first sorted the 11,390 deletions into $k = 8$, as well as $k = 16$ clusters using
 196 k-means clustering, based on frequency similarities across the 71 genomes (Supplemental Fig. 4, 5, 6A). We found
 197 that in some, albeit not in all deletion clusters, the Group 1 and 2 genomes visibly differ in their mean frequencies
 198 (Supplemental Fig. 4, 5). In total, 4 out of $k = 8$ clusters (comprising 14% of the total number of deletions) showed
 199 clear differences between Group 1 and Group 2 (effect size, Cohen's $d > 0.8$). We obtained similar results using for
 200 $k = 16$ (5 out of 16 clusters, comprising 11% deletions). We hypothesized that deletion loci in these 4 outlier clusters
 201 may have low mappability, i.e., the unique mapping potential of reads at a locus (Supplemental Fig. 7; Methods).
 202 Indeed, we found that the same clusters showing the strongest Group 1 vs. Group 2 differences also show the lowest
 203 average mappability (on average 16% and 24% lower mappability than rest of the clusters for $k = 8$ and for $k = 16$,
 204 respectively). We further reasoned that differential filtering of BAM data before publishing might explain different

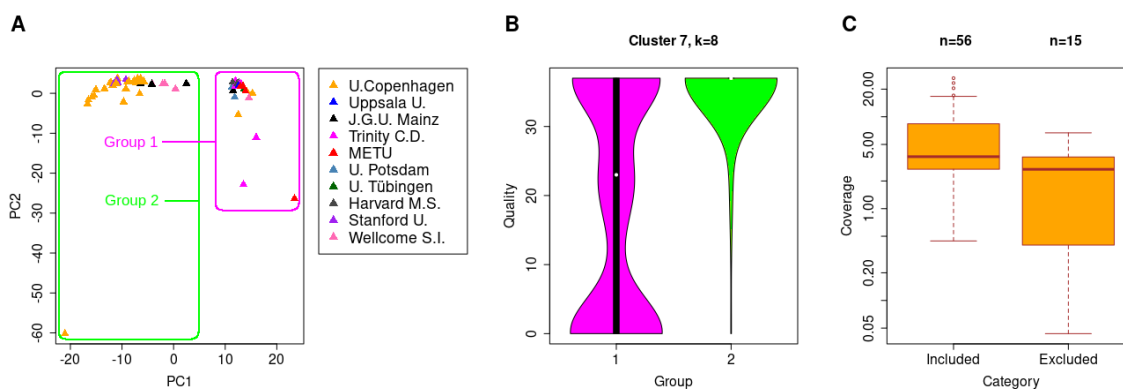


Figure 4: (A) A principal components analysis of deletion frequencies across the 71 genomes. There is significant distinction between the two groups observed in PC1 that explains 18% of the total variation (Wilcoxon rank sum test on PC1 values $p < 10^{-16}$). (B) Violin plot of MAPQ scores of 514 deletions belonging to outlier cluster 7 ($k = 8$), for Group 1 ($n = 30$) and Group 2 ($n = 41$) genomes. The two distributions are significantly different (Wilcoxon rank sum test, $p < 10^{-15}$). The same pattern is observed for the 2 other outlier clusters, which also happen to have low mappability, but not the remaining 5 clusters. (C) Boxplots representing the coverages of the included and excluded genomes in the refined dataset. The excluded genomes have significantly lower coverages on average (Wilcoxon rank sum test, $p = 0.01$).

205 deletion estimates at these loci between Group 1 vs. Group 2 genomes. Studying aligned read mapping quality
 206 (MAPQ) scores for reads mapping to each deletion in all 71 genomes, we discovered that Group 1 genomes had lower
 207 MAPQ scores than Group 2 genomes for deletions in the same 4 outlier clusters (Figure 4B; Supplemental Fig. 8, 9).
 208 This suggests that the observed laboratory-of-origin effect on deletion frequencies derives from differences among
 209 laboratories in BAM filtering before data submission: Group 1 data were not filtered for MAPQ, and consequently
 210 were genotyped accurately at low mappability deletions; hence their lower MAPQ scores at deletions in outlier clusters.
 211 In contrast, Group 2 BAM files were filtered for MAPQ, and at these loci, which lacked reads, CONGA genotyped
 212 artificial deletions.

213 If this assessment is true, the laboratory-of-origin effect could be removed by increasing the mappability score threshold
 214 per locus and filtering all BAM files for higher MAPQ. We thus increased the average mappability score threshold to
 215 ≥ 0.9 per deletion, thus discarding 12% of 11,390 deletions. We also filtered reads for $\text{MAPQ} \geq 30$. Using the 10,018
 216 deletions that remained after this double filtering, we repeated the principle components analysis (Supplemental Fig.
 217 6B) and replotted the heat map (Supplemental Fig. 3B). We now observed no laboratory-based grouping on the new PC1
 218 (Kruskal-Wallis test $p = 0.14$), and no difference in PC1 values between the previously identified Group 1 vs Group
 219 2 genomes (Wilcoxon rank sum test $p = 0.61$), suggesting that mappability and MAPQ filters effectively removed
 220 the laboratory-of-origin effect. We next tested each of the 10,018 deletions for a laboratory-of-origin effect on their
 221 frequencies; this identified only 130 (1.3%) showing a significant effect (Benjamini-Hochberg corrected Kruskal-Wallis
 222 test $p < 0.05$). We note that some of these cases could represent authentic population differences in deletion frequency,
 223 because data from different laboratories vary in their population origins.

224 To further ensure that the data is not influenced of major technical artifacts, we repeated the clustering analysis with the
 225 10,018 deletions that remained after filtering. This revealed that some singular ancient genomes behave as outliers in
 226 their deletion frequencies (Supplemental Fig. 4, 5). We could visually identify 15 such genomes out of the 71 total
 227 (Supplemental Table S2). A number of attributes could explain outlier behaviour. First, the coverage of the 15 outlier
 228 genomes was lower compared to the remaining 56 (Wilcoxon rank sum test, $p = 0.01$; Figure 6C). For instance, all
 229 three genomes with $< 0.1 \times$ coverage in our dataset (ne4, ko2, and DA379) were among the outliers. Second, the 15
 230 outlier genomes had on average shorter read length compared to the rest (median = 52.7 vs 65.8, Wilcoxon rank sum
 231 test, $p < 0.001$). One of these was the Iceman, with unusually short (50 bps) reads. Other characteristics could also
 232 play a role. Among the most extreme outliers was Bon002, which was the only sample produced using whole genome
 233 hybridization capture with MyBaits probes (Kılınc *et al.*, 2016), suggesting that the capture procedure may distort
 234 coverage. Consequently we decided to remove these 15 genomes from further analysis.

235 2.5 A comparison of deletion and SNP diversity across 56 ancient genomes

236 The above genotyping and filtering steps resulted in a refined set of 10,018 autosomal deletions ascertained in present-day
237 Africa and genotyped in 56 ancient Eurasian genomes, with 328-649 deletions (median = 411) detected in heterozygous
238 or homozygous state per genome, and 28% detected in at least one genome. We used this dataset to test two hypotheses
239 on deletion frequencies: (i) that deletion diversity patterns will broadly parallel SNP diversity patterns as previously
240 reported for modern-day CNV datasets (Conrad and Hurler, 2007; Levy-Sakin *et al.*, 2019; Almarri *et al.*, 2020), and
241 (ii) that deletion frequencies will reflect some degree of negative selection, which may be expected as deletions can bear
242 harmful due to various effects including expression level changes, exon loss, or frame-shifting.

243 To investigate the first hypothesis, we compared pairwise genetic distances among the 56 individuals (Figure 5A)
244 calculated using either SNPs or deletion genotypes. For this, we collected 5,991,735 autosomal SNPs ascertained
245 in African Yoruba individuals in the 1000 Genomes Dataset and genotyped our 56 ancient genomes at these loci
246 (Methods). We then calculated pairwise outgroup-f3 statistics, a measure of shared genetic drift between a pair of
247 genomes relative to an outgroup population (Patterson *et al.*, 2012). Outgroup-f3 values were calculated for all pairs of
248 ancient genomes using either SNPs or deletions, and these were then used to calculate SNP -or deletion- based
249 genetic distances between each pair (1-f3). We observed strong positive correlation between the two resulting distance
250 matrices (Spearman $r = 0.71$, Mantel test $p = 0.001$) (Figure 5B). We further investigated relationships among ancient
251 individuals by summarizing these distances using multidimensional scaling (MDS). This revealed highly similar patterns
252 in the SNP -and deletion- based MDS plots, with clear clustering among west and east Eurasian genomes (Figure 5C,
253 D). These results support the notion that deletion diversity patterns are at least partly shaped by demographic history
254 and shared genetic drift.

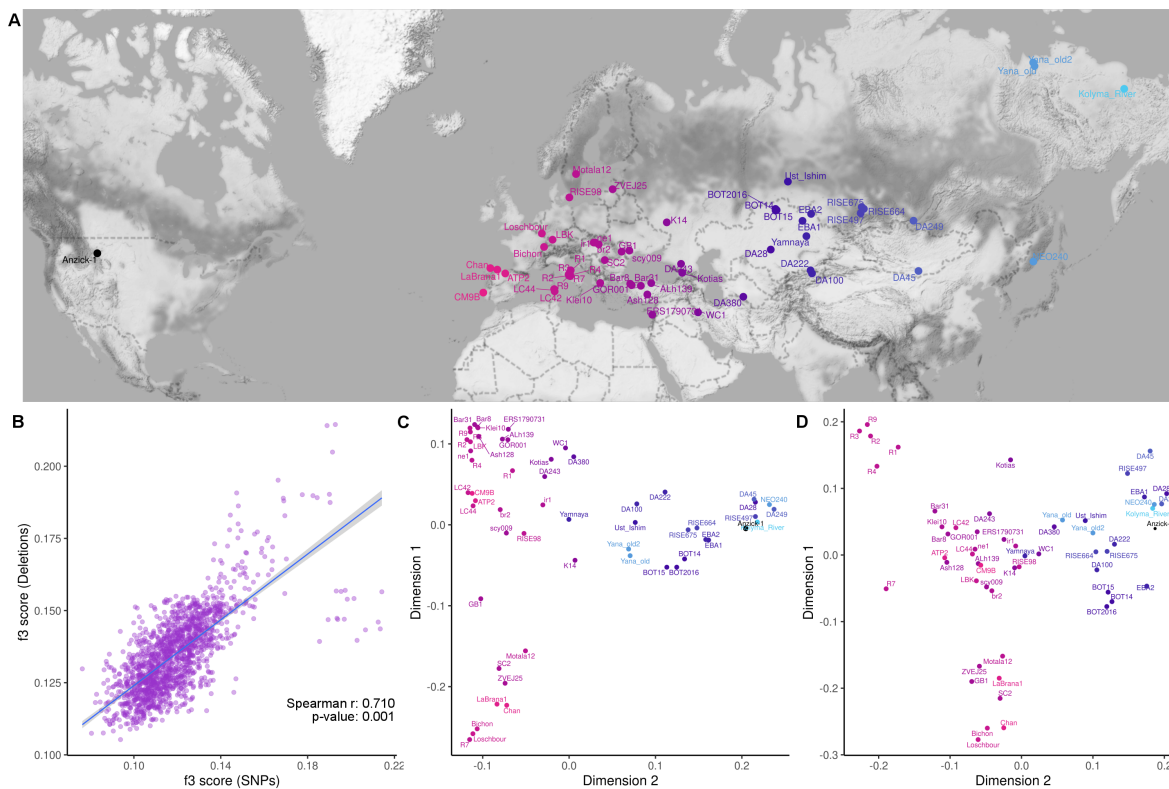


Figure 5: (A) Geographic locations of the 56 ancient individuals. (B) Comparison of genetic distances calculated using SNPs and deletions. We calculated the Spearman correlation coefficient between two matrices and then calculated Mantel test p-value using the "mantel" function in R package "vegan" (v2.5-7). (C) and (D) represent multidimensional scaling plots that summarize outgroup-f3 statistics calculated across all pairs among the 56 ancient individuals using SNPs and deletions, respectively.

255 2.6 Negative selection on deletion variants

256 We next studied the site-frequency-spectrum (SFS) of autosomal deletions across these 56 genomes in comparison
257 with the SNP SFS. For this, we first polarised deletions as well as SNPs using chimpanzee and bonobo genomes
258 to represent the ancestral state, thus determining 9,579 deletions and 4,907,535 SNPs derived in the human lineage
259 (Methods). Second, to allow comparison with the pseudo-haploidized SNP genotype data, we randomly chose one allele
260 per genome (i.e., deletion or no event) in the deletion dataset. Set side by side with the SNP SFS, we observed a stark
261 excess of singletons among deletions, which is consistent with stronger negative selection on the latter (Figure 6A).

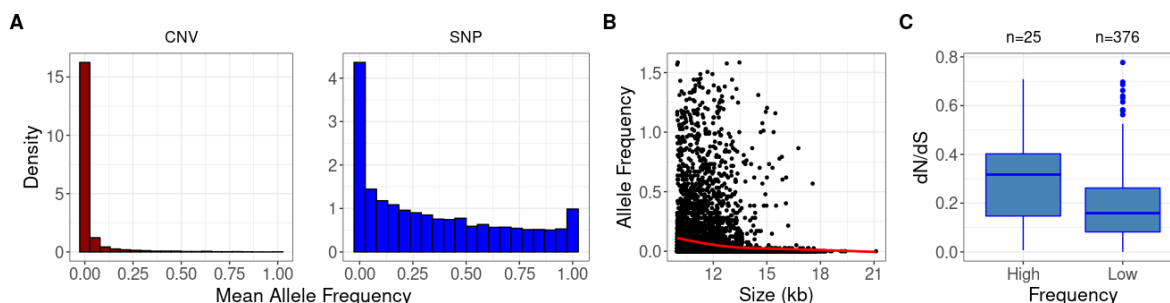


Figure 6: (A) The site-frequency-spectra of derived deletion alleles (on the left) and derived SNP alleles (on the right). The x-axes show mean allele frequency for each locus calculated using only those genomes where a locus has been observed (e.g. an allele observed in 10 out of 40 genomes will be represented as 25%). The two distributions are significantly different from each other (Kolmogorov-Smirnov test $p < 10^{-15}$). (B) The size distribution (in Kbps) of the deletions versus mean allele frequency. The red line shows the fitting of smoothing spline and indicates a negative correlation (Spearman correlation $r = -0.35$, $p < 10^{-16}$). Both axes were \log_2 -scaled. (C) Boxplots representing the dN/dS values of high ($n = 25$) vs. low ($n = 376$) frequency deletions.

262 If deletions are indeed under strong negative selection, as indicated by their SFS, we may expect longer deletions, as
263 well as deletions containing evolutionary conserved genes, to be segregating at lower frequencies. To test the first
264 idea, we compared deletion allele frequencies with their length. As hypothesized, deletion size and its frequency were
265 negatively correlated across the 56 genomes (Spearman correlation $r = -0.35$, $p < 10^{-16}$) (Figure 6B).

266 To test the second hypothesis, we first determined annotated genes overlapping with our deletion dataset and found
267 461 Ensembl (v75) human genes. Overall, 5% of the 9579 derived deletions overlapped with at least one gene, and
268 among those deletions, each overlapped with an average of 1.33 genes. We then collected mouse-human dN/dS ratios
269 (Methods) for these genes ($n = 401$, 0-0.77, median = 0.16, mean = 0.20), which is an inverse measure of protein
270 sequence conservation. We tested for a difference between lower frequency (below average) and higher frequency
271 (above average) deletions in terms of dN/dS values of the genes they overlapped with. We found that deletions with
272 lower allele frequency also had lower dN/dS values (Mann-Whitney U test, two-sided $p = 0.002$; Figure 6C), with an
273 effect size of Cohen's $d = 0.68$. This suggests that deletions disrupting highly conserved genes tend to be rarer, again
274 consistent with the action of negative selective.

275 2.7 Time and memory consumption

276 Finally, we examined time and memory requirements of CONGA. We first tested our performance of deletions with
277 BAM files of the 71 ancient genomes presented above. This finished in ~ 12 hours in total with as low as 2.2 GB of
278 peak-memory consumption. This is ~ 10 minutes per genome. In order to evaluate CONGA's performance with a
279 higher coverage genome sample, we ran 30 genomes (randomly selected 10 CEU, 10 YRI, 10 TSI) from the 1000
280 Genomes Project Phase 3, which had mean $7.4\times$ coverage (Sudmant *et al.*, 2015). The analysis took just slightly longer,
281 ~ 15 minutes average per genome, with similar memory usage.

282 We also compare the time and memory requirements of CONGA, FREEC and CNVnator in Table 2. In order to
283 benchmark these tools, we used a $5\times$ simulated genome (the same genome with medium sized CNVs used in the
284 simulation experiments described above) with the same computing resources¹. CONGA was $2\times$ slower than the other
285 two software in time performance when run in default mode, i.e. when genotyping duplications and deletions together.
286 The reason for this is CONGA's use of its own small-scale read mapper for split-reads, which creates the bottleneck for
287 time and memory usage. However, when only genotyping deletions, split-read information is not utilized; in this case
288 CONGA's running time reduces to $\sim 1/3$ and memory usage to $\sim 1/7$ of the other two algorithms.

¹Intel(R) Xeon(R) CPU E5-2640 v2 @ 2.00GHz: 2CPUs * 8 cores each=16 cores total and 216GB RAM

289 We further provide a comparison of CONGA's performance on genomes of various depths of coverage in Supplemental
290 Table S1.D, calculated using the down-sampled $23\times$ Yamnaya genome (with coverages between $23\times$ and $0.07\times$).

Table 2: Time and Memory Consumption

Tools	Time (h:m)	Peak Memory Usage (GB)
CNVnator	0:32	12
FREEC	0:39	7.1
CONGA (deletions only)	0:12	1.2
CONGA (deletions and duplications)	1:19	10

Time and memory consumption of each algorithm for a simulated genome of $5\times$ depth of coverage. "Time" refers to wall clock time and "Peak memory usage" is the maximum resident set size.

291 3 Discussion

292 Modern human genome sequencing experiments today typically reach coverages $>20\times$ and increasingly use long read
293 technology, and such experiments can employ diverse read signatures to reliably identify CNVs (Alkan *et al.*, 2011).
294 CONGA's approach that mainly relies on the read-depth signature is naive in comparison; however, using read-depth
295 appears as the main practical solution given the short fragment size and the predominance of low coverage (around or
296 $<1\times$) among ancient genome datasets.

297 3.1 CONGA's overall performance and utility

298 Despite these challenges, our experiments using simulated genomes and down-sampled real ancient FASTQ data
299 showed that CONGA can relatively efficiently genotype deletions and duplications of size >1 Kbps at $1\times$ coverage,
300 or even lower. CONGA outperformed two "modern DNA" CNV discovery algorithms, FREEC and CNVnator, two
301 methods previously employed in ancient genome analyses (Smith *et al.*, 2017; Bhattacharya *et al.*, 2018). CONGA
302 exceeded both tools in true positive as well as true negative rates, especially at coverages $<1\times$. This is unsurprising, as
303 these latter tools were developed for the purpose of discovering novel CNVs in relatively high coverage genome data.
304 In contrast, CONGA has been specifically developed for analyzing low coverage genomes by genotyping CNVs at
305 previously ascertained loci.

306 CONGA further achieved acceptable accuracy ($\sim 50\%$ TPR and $<30\%$ FDR) in deletion copy number estimates in
307 genomes of $0.5\times$ coverage. At lower depths of coverage and also when genotyping deletions <1 Kbps, recall and/or
308 precision were weaker. CONGA's performance on duplications was also weaker, as we discuss below.

309 Nevertheless, the relatively high accuracy at $0.5\times$ coverage suggests that the method could be used to genotype deletions
310 across a considerable fraction of published shotgun sequenced ancient genomes. Beyond aDNA, CONGA will also be
311 suitable for CNV analyses for any low depth whole-genome sequencing (WGS) experiment. Such studies are increasing
312 in number due to the trade-off between budget limitations and the wealth of genome-wide information that can be used
313 in population and conservation genetics (e.g. (Vieira *et al.*, 2016)).

314 3.2 Caveats in duplication genotyping

315 In simulated genome experiments, CONGA's performance in genotyping duplications was similar to that in deletions.
316 Here, read-depth information provided comparable power as in genotyping deletions; furthermore, paired-end informa-
317 tion from split reads also yielded results with acceptable accuracy (c.50%), although only at $5\times$ coverage (Supplemental
318 Table S1.E). In down-sampling experiments, CONGA showed slightly lower performance in duplication genotyping than
319 in deletion genotyping when using two ancient genomes available as FASTQ files. However, CONGA's performance
320 was dramatically low on the $23\times$ ancient BAM file, Yamnaya. Our results suggest the following reasons behind this
321 discrepancy: (i) The Yamnaya genome was available as a BAM file, and we presume it was subject to pre-publication
322 processing (at the alignment or quality filtering steps) that removed excess reads at duplicated loci. (ii) Consequently,
323 97% of duplications CONGA genotyped in the original ($23\times$) BAM file were called using paired-end information,
324 instead of read-depth information. (iii) Because paired-end information is more severely affected than read-depth
325 information with decreasing coverage (as the former uses information from only reads that overlap breakpoints, while
326 the latter uses information from across the length of the CNV), genotyping duplications in this BAM files became
327 infeasible at $<5\times$ coverage.

328 The majority of shotgun ancient genomes in public databases are only published as BAM files. The motivation is
329 to avoid publishing environmental DNA reads, which constitute a large fraction of shotgun sequencing experiments.

330 The majority of published files are also at $<5\times$ coverage. Hence, most published ancient shotgun genomes are not
331 amenable to duplication genotyping with CONGA. This is highly unfortunate, as gene duplications are a major source
332 of evolutionary adaptation that would be valuable to study also in ancient populations (see below).

333 3.3 Caveats in deletion genotyping

334 Applying CONGA to genotype deletions on 71 real ancient shotgun genomes from North Eurasia and the Americas
335 using a dataset of modern-day African deletions yielded further critical observations. First, we found that genotyping
336 deletions on published BAM files can also suffer from artifacts caused by pre-filtering of published reads based on
337 mapping quality, a practice which apparently varies among research groups. Pre-filtered BAM files contain a paucity of
338 reads originating from loci with relatively low mappability, creating artificial deletion signals. Nevertheless, we find
339 that this issue can be overcome by increasing the average mappability threshold for deletion loci in the reference set and
340 by applying higher mapping quality filters to all BAM files - a solution not available for genotyping duplications.

341 Another observation was that 15 of the 71 genomes analysed displayed unexpectedly high deletion frequencies at
342 loci that are rarely or never genotyped in the rest of the 71 individuals. Given the close evolutionary relationship
343 among Eurasian human populations, the majority of these signals are unlikely to be authentic, but rather originate from
344 experimental artifacts and/or variability of DNA preservation among samples. Indeed, we could notice particularities in
345 experimental protocols for some of these genomes, such as lower coverage, shorter read lengths, or the application of
346 whole-genome hybridization capture.

347 These observations on real genomes indicate that $0.1\times$ coverage may be the lower threshold for deletion genotyping of
348 >1 Kbps events, which is also consistent with the simulation results. Whole-genome hybridization capture may also be
349 incompatible with reliable CNV detection. That said, we still lack clear explanations for outlier deletion frequency
350 patterns for some of these 15 genomes. For instance, the genome SI-45 has coverage $>3\times$ and an average read length
351 of 60 bps, but nevertheless displays unusual deletion patterns. We suspect that such unexpected patterns might reflect
352 technical peculiarities in DNA extraction, library preparation, sequencing or data filtering.

353 3.4 Community recommendations for improving CNV analyses in ancient genomes

354 We mark the urgent need for changes in practice in producing and publishing ancient genomes to allow reliable study of
355 CNVs in addition to SNPs.

- 356 • Most published ancient genome data to date is SNP capture data, which is largely worthless for CNV analyses.
357 Our results underscore the long-term value of shotgun sequencing data over SNP capture data. We also note
358 that whole-genome capture data may be incompatible with CNV discovery.
- 359 • Researchers should strive to publish their data either as FASTQ files, or as BAM files aligned with the most
360 relaxed parameters and without applying any filtering. Our results show that this is absolutely necessary to
361 allow duplication genotyping at low coverage and to avoid biases in deletion genotyping on BAM files.
- 362 • Sharing all details on DNA extraction, library construction, as well as the alignment and preprocessing steps
363 used in creating the exact version of datasets submitted to public databases is crucial for healthy reuse of the
364 data.

365 We further recommend the following for comparative CNV analyses in ancient genomes:

- 366 • Studying and –if necessary– normalizing mapping quality distributions across published BAM files is crucial
367 to avoid biases.
- 368 • Given the unavoidable heterogeneity of ancient genome data sources and the sensitivity of CNV detection and
369 genotyping tools to such heterogeneity, researchers should pay utmost care to identify and filter out abnormal
370 patterns of CNV frequencies in joint datasets.

371 3.5 Studying drift and selection in past populations

372 Beyond these technical aspects, our analysis of >1 Kbps deletions genotyped in 56 ancient genomes revealed how
373 this CNV species has been evolving under genetic drift and negative selection, simultaneously. First, we found that
374 genetic distances based on deletion frequencies were overall strongly correlated with those of SNP frequencies, and that
375 they reflected spatial distances among populations. This is consistent with the observation that CNV frequencies in
376 modern-day human populations at least partly reflect demographic processes, such as geographical isolation, drift and
377 admixture (Conrad and Hurler, 2007; Levy-Sakin *et al.*, 2019; Almarri *et al.*, 2020).

378 Second, we observed that the deletion frequency spectrum in ancient Eurasians was considerably steeper than that
379 of SNPs. While a notable proportion of derived SNPs polymorphic in Africans were fixed in our ancient Eurasian
380 sample, we found no such fixed deletion event. We also found that large deletions, as well as deletions disrupting
381 highly conserved genes, segregated at significantly lower frequencies compared to the rest of the deletion set. These
382 results indicate that deletions are on average less evolutionarily neutral than SNPs and evolve under negative selection
383 pressure. Intriguingly, these deletions have apparently continued to segregate across continents for possibly more than
384 hundred thousand years, perhaps aided by bottleneck events. CONGA thus opens up a new avenue of research into
385 deleterious CNV mutation loads in natural populations, which can now be studied using the rapidly growing number of
386 low coverage WGS experiments.

387 4 Methods

388 Among various approaches developed for CNV discovery using high throughput sequencing data, almost all use the
389 fact that read-depth, i.e., the density of reads mapped to the reference genome, will be on average lower in deleted
390 regions and higher in duplicated regions (Alkan *et al.*, 2011; Ho *et al.*, 2020). The distance between paired-end reads,
391 their orientation, and split-read information (start and end of reads mapping to different locations) are further sources
392 of information used in determining CNVs. Although available CNV discovery algorithms generally perform well in
393 modern-day human genome sequencing data with high coverage, this is not necessarily the case for ancient genomes, as
394 well as other low coverage sequencing experiments (Supplemental Fig. 10, 11). The first reason is that the majority of
395 shotgun ancient genomes are produced at low coverage (typically $<1\times$), which limits the use of read-depth information.
396 Second, ancient DNA fragments are short and of variable size (typically between 50-100 bps) (Shapiro and Hofreiter,
397 2014). Thus, paired-end information is absent, and available split-read information is also limited. Variability in ancient
398 DNA preservation and genome coverage (Pedersen *et al.*, 2014) is yet another noise source that is expected to limit
399 efficient CNV discovery. CONGA overcomes these limitations using genotyping instead of *de novo* discovery. It
400 estimates whether a candidate CNV, the location of which is provided as input, is present in a genome in BAM format. It
401 also estimates the genotype, i.e., the heterozygous or homozygous state. CONGA makes use of read-depth information
402 for deletions, and both read-depth and split-read information for duplications.

403 4.1 Likelihood-based read-depth calculation for deletion and duplication genotyping

404 The input to the algorithm is (1) a list of candidate CNV locations and CNV type, i.e., deletion or duplication, and (2) a
405 data set of reads aligned to the linear reference genome, e.g., using BWA (Li and Durbin, 2009), which should be in
406 BAM format.

407 In order to calculate the likelihood of a CNV at a given locus based on read-depth information, CONGA uses an
408 approach akin to (Soylev *et al.*, 2019). Let (S_i) be the i^{th} CNV in our CNV input list, defined by the breakpoint interval
409 (B_l, B_r) and the type of CNV: a deletion or duplication. At this locus, CONGA calculates the likelihood of the three
410 possible genotype states, k , given the read alignment data and CNV type. The genotype states are: no event ($k = 0$), a
411 heterozygous state ($k = 1$), or a homozygous state ($k = 2$). The likelihood, in turn, is calculated by comparing the
412 observed (O_i) read-depth versus the expected (E_{ik}) read-depth within (B_l, B_r), given the three different genotypes.
413 We detail the steps below.

- 414 1. We count the total number of mapped reads within that locus (falling fully within the interval (B_l, B_r)). This
415 is the observed read-depth, (O_{RD}).
- 416 2. We calculate expected read-depth under a no event scenario, i.e., representing the diploid state. Here we
417 account for the GC bias in high throughput sequencing data (Smith *et al.*, 2008), by using LOESS smoothing
418 to normalize read-depth for GC content. Specifically, for each chromosome, we calculate the read-depth values
419 per GC percentile for sliding windows of size 1,000 bps (step size = 1 bp). We then calculate the average
420 read-depth per GC percentile. Then, using the chromosome-wide average GC value for the interval (B_l, B_r),
421 we calculate the expected diploid read-depth, $E_{ik=0}$.
- 422 3. We model the read-depth distribution as Poisson, using the expected read-depth values for $k = 0$, $k = 1$,
423 $k = 2$. We calculate the probability $P(RD_{S_i} | state = k)$ as:

$$P(RD_{S_i} | state = k) = \frac{E_{ik}^{O_i} \times e^{-E_{ik}}}{O_i!},$$

424 where E_{ik} is the expected read-depth given $state = k$, and O_i is the observed read-depth at that specific locus.
425 A typical autosomal human locus is diploid (has copy number = 2); therefore when there is no CNV event
426 ($k = 0$), the expected value of O_i should be $E_{ik=0}$.

427 If a genome is homozygous for a deletion, we expect no reads mapping to the region, thus $O_i \sim E_{i_{k=2}} = 0$. For
428 heterozygous deletions, the expected number of mapped reads in that interval will be half of the expected diploid
429 read-depth: $O_i \sim E_{i_{k=1}} = E_{i_{k=0}}/2$. For homozygous duplications, we expect $O_i \sim E_{i_{k=2}} = E_{i_{k=0}} \times 2$. For
430 heterozygous duplications, we expect $O_i \sim E_{i_{k=1}} = E_{i_{k=0}} \times 1.5$.

431 4. We calculate a likelihood-based score, which we term the C-score, to estimate how likely locus S_i carries a
432 non-reference variant in a genome, in either one copy or two copies. For this we use the calculated likelihoods
433 for the three states. We define the C-score as the maximum of the likelihoods of (S_i) being present in
434 heterozygous state ($k = 1$) or in homozygous state ($k = 2$) in that genome, over the likelihood of no event
435 ($k = 0$). We use the log function to avoid numerical errors.

$$C - score(S_i) = \frac{\max(\log(P(RD_{S_i}|k=1)), \log(P(RD_{S_i}|k=2)))}{\log(P(RD_{S_i}|k=0))},$$

436 The C-score is distributed between 0 and $+\infty$, with lower scores indicating higher likelihood of a true CNV
437 event.

438 Alternative approaches could also be considered within the same framework. One is using the negative binomial to
439 model the read-depth (Miller *et al.*, 2011). However, the Poisson is simple and a most commonly used tool, and
440 our simulations (described in the following) empirically demonstrate its effectiveness. Another approach could be
441 calculating the likelihood of >2 copies of duplication events, such as multicopy genes (Sudmant *et al.*, 2010). These
442 alternatives may be considered in future work.

443 4.2 Split-read and paired-end signatures for duplication genotyping

444 Beyond read-depth, information of paired-end reads or read fragments that do not linearly map to the genome can be
445 used to identify CNVs. Ancient genomes are sometimes single-end and sometimes paired-end sequenced, but in the
446 latter case, short overlapping reads are typically merged into a single read before alignment. Ancient genome data is
447 thus practically single-read. However, the split-read method can be applied on single-read ancient genome data, which
448 emulates paired-end information for genotyping duplications. This approach is visualized in Figure 7. We therefore
449 designed CONGA to include both paired-end and single-end reads as input and evaluate the paired-end signature
450 information.

451 First, assume a read of length L mapped to position pos_x in the reference genome, where pos_x is assumed to be one
452 of the breakpoints of a putative CNV. There always exists a subsequence $\geq L/2$ that will have at least one mapping
453 in the reference genome with some error threshold. Thus, we can split a read into two subsequences, assigning the
454 actual mapping to one of the pairs and remapping the other subsequence ("split segment") as a second pair. There
455 are two possible split strategies: an even decomposition, where both subsequences are of equal lengths, or an uneven
456 decomposition, where the subsequences are of unequal lengths. Given the infeasibility of testing each split position
457 and the fact that ancient reads are typically already short, we follow (Karakoc *et al.*, 2012) and split the read from the
458 middle to obtain two reads with equal lengths $L/2$. If a read overlaps a duplication breakpoint, and assuming that the
459 expected position of the breakpoint will be uniformly distributed within the read, the split segment will map to the
460 reference genome with insert size—the distance between the split-read pairs—greater than zero.

461 With this simple observation, the need to observe all possible breakpoints can be eliminated. Thus, given a single-end
462 read Rse_i , we define $Rpe_i = (l(Rpe_i[pos_x : pos_x + RL/2])$ and $r(Rpe_i[pos_y : pos_y + RL/2]))$, where pos_x is the
463 initial mapping position of the single-end read, pos_y is the remapping position of the split read, RL is the length of the
464 single-end read observed before the split, $l(Rpe_i[pos_x : pos_x + RL/2])$ is the left pair within pos_x and $pos_x + RL/2$
465 and $r(Rpe_i[pos_y : pos_y + RL/2])$ is the right pair within pos_y and $pos_y + RL/2$ of the paired-end reads. We use this
466 information as described in the following section.

467 4.2.1 Remapping paired-reads and utilizing paired-read information

468 According to our remapping strategy, we use a seed-and-extend approach similar to that implemented in mFAST (Alkan
469 *et al.*, 2009), where a read is allowed to be mapped to multiple positions. Our main concern here is that the split
470 segment, due to its short length, can be mapped to unrealistically high numbers of positions across the genome. To
471 overcome this problem we use the approach developed in TARDIS (Soylev *et al.*, 2017), allowing the split segment to
472 be mapped only up to 10 positions within close proximity (15 Kbps by default) of the original mapping position and
473 applying a Hamming distance threshold for mismatches (5% of the read length by default).

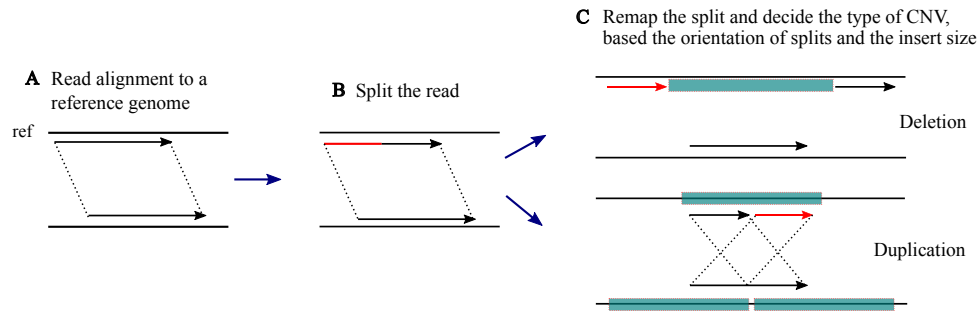


Figure 7: The figure shows our split read approach to emulate paired-end using single-end reads. (A) We are given short-read Illumina mappings in a BAM file as input. The top horizontal line represents the reference genome sequence ("ref"), while the lower line represents the donor genome. (B) We split the read from the middle, keeping the initial mapping as one element and the other subsequence (split segment) as the second element of a pair. (C) We remap the split segment to the reference genome, and evaluate the position and the orientation of both reads to identify the presence of putative CNVs.

474 Based on the distance between the reads (insert-size) and orientation, we then evaluate the type of putative CNV.
475 As Figure 7C shows, if the split segment maps behind the initially mapped segment of the same pair to generate a
476 reverse-forward mapping orientation, this would be an indication of a duplication.

477 In order to utilize this paired-read information, for each CNV locus used as input to our algorithm, we count the number
478 of read-pair (i.e. split segments) that map around ± 5 Kbps of the breakpoints. Each such read-pair is treated as one
479 observation. We use these counts in combination with the C-score (read-depth information) to genotype duplications
480 (see below). We do not use this read-pair information for genotyping deletions due to its low effectiveness in our initial
481 trials (Supplemental Table S1.E).

482 4.3 Mappability filtering

483 The probability of unique alignment of a read of certain size varies across the genome, mainly due to repetitive
484 sequences. Various algorithms estimate this probability, termed mappability, across the genome for k-mers of specific
485 length (Koehler *et al.*, 2011; Derrien *et al.*, 2012; Karimzadeh *et al.*, 2018; Pockrandt *et al.*, 2020). This is calculated
486 by extracting k-mers of given length through the genome, remapping them to the reference genome, and measuring
487 mappability as the proportion of unique mappings (Karimzadeh *et al.*, 2018). Because low mappability regions can
488 be confounded with real deletions, we use mappability information to filter out CNV loci that could represent false
489 positives.

490 CONGA accepts any mappability file in BED format, where values are distributed between 0 and 1. These can then be
491 used to filter out CNVs for minimum mappability.

492 In our experiments, we used the 100-mer mappability data from the ENCODE Project (ENCODE Project Consortium,
493 2012), available at <http://hgdownload.cse.ucsc.edu/goldenpath/hg19/encodeDCC/wgEncodeMapability/>. Using this
494 data, for each CNV event (S_i), we calculated the average mappability value within its breakpoints. We used either of
495 two different minimum average mappability thresholds for the CNV events we analyzed: 0.5 (relaxed) or 0.9 (strict).
496 Our deletion frequency analysis results show that the strict filter is more reliable especially when analyzing data sets of
497 heterogeneous origin (see Results).

498 4.4 Simulation and down-sampling experiments

499 4.4.1 Simulating ancient genomes with implanted deletions and duplications

500 Our goal here was to study the performance of CONGA on different sized deletions or duplications using simulated
501 genomes containing implanted CNVs and to determine thresholds for reliably calling these variants. We first employed
502 VarSim (Mu *et al.*, 2015) to simulate and insert deletions and duplications into the human reference genome GRCh37.
503 We repeated this three times, for small (100 bps - 1000 bps), medium (1000 bps - 10,000 bps), and large (10,000 bps -
504 100,000 bps) CNVs. As a result we generated three CNV implanted genomes, with around 1500 deletions and 1500
505 duplications each (between 1385 and 1810). The CNVs were produced so that they were non-overlapping, and their
506 length distribution and exact counts are provided in Supplemental Fig. 1.

507 To evaluate specificity in addition to sensitivity, we also included a background (false) CNV set in the experiment,
508 which would not be implanted but would be queried as part of the candidate list. This background set was prepared
509 using recently published deletion and duplication calls from human genome sequencing experiments (Audano *et al.*,
510 2019; Chaisson *et al.*, 2019; Zook *et al.*, 2020; Collins *et al.*, 2020) and also sequencing data from African populations
511 (AFR) from Phase 3 of the 1000 Genomes Project (Sudmant *et al.*, 2015). We compiled a list of 17,392 deletions
512 and 14,888 duplications that were non-overlapping and of size $> \sim 1000$ bps using BEDTools mergeBed (Quinlan and
513 Hall, 2010). When evaluating genomes with small CNVs (100 bps - 1,000 bps), we additionally included small CNVs
514 from (Chaisson *et al.*, 2019). Specifically we added 4,623 deletions and 3,750 duplications of size 100 bps - 1,000 bps
515 to the above background list.

516 In order to assess CONGA's performance, we added the true CNVs generated using VarSim to this background set,
517 such that only $\sim 10\%$ of the input candidate CNV list were true events. Finally, we determined how many of these true
518 events could be correctly called by CONGA and other software.

519 4.4.2 Simulating ancient genome read data

520 We then used the above-described simulated genomes as input to Gargammel (Renaud *et al.*, 2017), which generates
521 ancient-like Illumina reads, i.e., short reads of variable size bearing postmortem damage (i.e., C-to-T transitions at
522 read ends) and including adapters. As Gargammel is able to generate aDNA fragments with size distribution given as
523 input, we used a subset of (Fu *et al.*, 2014), which is its default for this software. We used Gargammel to produce reads
524 at various depths of coverage: $0.05\times$, $0.1\times$, $0.5\times$, $1\times$ and $5\times$. We then removed adapters and merged overlapping
525 reads (Schubert *et al.*, 2016) to generate single-end Illumina reads. These reads had sizes ranging between 34 bps
526 and 139 bps, with average 69 bps and median 66 bps (calculated using $1\times$ coverage data, but others also had similar
527 distributions). We mapped the Gargammel-output reads back to the human reference genome (hg19, or GRCh37) using
528 BWA-aln (Li and Durbin, 2009) with parameters "-l 16500 -n 0.01 -o 2" (Supplemental Material). Note that BWA-aln
529 has been shown to be more accurate for short ancient reads than BWA-mem (Oliva *et al.*, 2021).

530 4.4.3 Evaluation of CONGA, CNVnator and FREEC with simulated ancient genome data

531 We ran CNVnator (Abyzov *et al.*, 2011) and FREEC (Boeva *et al.*, 2012) on the simulated genomes with parameters
532 described in the Supplemental Material and CONGA with two values for the C-score (< 0.3 and < 0.5). We used the
533 above-described list of CNVs as the input candidate set.

534 To determine true calls, we used $> 50\%$ reciprocal overlap for the two CNV events (the event in the input event set and
535 the called event) to be considered the same. This calculation was done using BEDTools (Quinlan and Hall, 2010). Note
536 that the number of True CNVs are: 1810 deletions and 1751 duplications for 100 bps - 1000 bps; 1680 deletions and
537 1684 duplications for 1000 bps - 10,000 bps; and 1385 deletions and 1532 duplications for 10,000 bps - 100,000 bps.

538 4.4.4 Down-sampling experiment with real ancient genomes

539 We used three relatively high coverage ($\sim 23.3\times$, $\sim 13.1\times$ and $\sim 9.6\times$ respectively) genomes of a Yamnaya culture-
540 related individual from early Bronze Age Karagash (hereafter Yamnaya), Kazakhstan (de Barros Damgaard *et al.*,
541 2018b), a Saqqaq culture-related individual from Bronze Age Greenland (hereafter Saqqaq) (Rasmussen *et al.*, 2010),
542 and a 4500-year old East African hunter-gatherer individual from Mota Cave in Ethiopia (hereafter Mota) (Llorente
543 *et al.*, 2015). Using this data, and the above-described 17,392 deletions and 14,888 duplications of size > 1 Kbps (see
544 above) as input, we genotyped 2639 deletions and 1972 duplications in Yamnaya (deletion sizes: 1 Kbps to 4 Mbps,
545 median = 4 Kbps, mean = 23 Kbps; duplication sizes: 1 Kbps to 28 Mbps, median = 14 Kbps, mean = 80 Kbps); 1581
546 deletions and 4097 duplications in Saqqaq (deletion sizes: 1 Kbps to 5 Mbps, median = 5 Kbps, mean = 17 Kbps;
547 duplication sizes: 1 Kbps to 28 Mbps, median = 16 Kbps, mean = 70 Kbps); and 688 deletions and 638 duplications in
548 Mota (deletion sizes: 1 Kbps to 130 Kbps, median = 4 Kbps, mean = 7 Kbps; duplication sizes: 1 Kbps to 28 Mbps,
549 median = 6 Kbps, mean = 82 Kbps).

550 We then randomly down-sampled the BAM files to various depths using Picard Tools (Pic, 2019): between $16-0.07\times$
551 for Yamnaya; $9-0.05\times$ for Saqqaq; $7-0.03\times$ for Mota. We note that this down-sampling procedure does not produce the
552 exact targeted depths, which is the reason why we have variable coverages in Fig. 3.

553 For calling deletions we used C-score < 0.5 . For calling duplications, we called events that fulfilled either of the
554 following conditions (a) C-score < 0.5 , or (b) C-score < 10 and read-pair support > 10 . Finally, treating the results of
555 the original data as the correct call-set, we calculated TPR (true positive rate) and FDR (false discovery rate) for
556 the down-sampled genomes. We considered CNVs with $\geq 50\%$ reciprocal overlap as representing the same event,
557 calculated using BEDTools (Quinlan and Hall, 2010).

558 **4.4.5 C-score and read-pair cutoffs and minimum CNV size**

559 We ran CONGA with a range of parameter values for the C-score [0.1-10] and for minimum read-pair support (from 0
560 support to >30), and using the above-described true event sets as the input candidate set involving medium and large
561 CNVs (1680 deletions and 1684 duplications for 1000 bps - 10,000 bps, and 1385 deletions and 1532 duplications for
562 10,000 bps - 100,000 bps).

563 We used simulation results (Supplemental Table S1.E) to choose an effective an cutoff for calling CNVs. For both
564 deletions and duplications, we decided to use C-score <0.5, which appears to yield a good trade-off between recall and
565 precision. Specifically, in simulations, this cutoff ensured an F-score of >0.5 at 0.1× for >1 Kbps deletions, and much
566 higher F-scores at higher coverages (Supplemental Fig. 12).

567 In addition, we observed that read-pair support >10 could be useful for identifying duplications in the absence of
568 read-depth support, but only when coverages were $\geq 1 \times$ (Supplemental Table S1.E; Supplemental Fig. 13). Moreover,
569 read-pair support was not valid for detecting deletions.

570 We note that CONGA outputs the C-scores and read-pair counts for all input CNVs. Users can choose alternative
571 cutoffs to increase recall (higher C-scores) or precision (lower C-scores).

572 The simulation experiments showed that CONGA was not efficient in identifying events <1 Kbps. We therefore
573 designed CONGA not to evaluate events <1 Kbps under default parameters. This can be modified by the user if needed.

574 **4.5 Analysis of real ancient genomes**

575 **4.5.1 Ancient genome selection and preprocessing**

576 We selected 71 ancient shotgun or whole-genome captured genomes from individuals excavated in West and East
577 Eurasia and in North America (Supplemental Table S2). Our sample set belongs to a time range between c.2,800-45,000
578 years Before Present (BP). Samples from 10 different laboratories were selected in order to study the effects of different
579 data production protocols on deletion genotyping. We also chose genomes with a range of coverage levels (0.04×-26×,
580 median = 3.45×) and that included both UDG-treated and non-UDG-treated libraries.

581 Selected ancient genomes were mapped to the human reference genome (hg19, or GRCh37) using BWA aln/samse
582 (0.7.15) (Li and Durbin, 2009) with parameters "-n 0.01, -o 2". PCR duplicates were removed using FilterUniqueSAM-
583 Cons.py (Kircher, 2012). We also removed reads with >10% mismatches to the reference genome and those of size
584 <35 bps. For preparing the "refined data set", reads with <30 mapping quality (MAPQ) were additionally removed.

585 **4.5.2 Candidate CNV call set for real ancient genomes**

586 Here our goal was to study properties of deletion variants in ancient genomes and to compare these with SNP variation
587 in terms of demographic history and purifying selection. Polymorphism data sets can suffer from ascertainment bias in
588 downstream evolutionary analyses (Clark *et al.*, 2005). A common practice to avoid this bias is to use SNPs ascertained
589 in a population that is an outgroup to the focal populations. We therefore used variants ascertained in modern-day
590 African populations for both calling SNP and deletion variants in our ancient genomes.

591 In order to create a candidate deletion call set to be used as input to CONGA, we downloaded deletions of size >1000
592 bps identified among 661 African population (AFR) genomes of the 1000 Genomes Project Phase 3 (Sudmant *et al.*,
593 2015). When a deletion was located inside the breakpoints of another deletion, we removed the internal one. In addition,
594 for pairs of deletions that had >50% overlap, we filtered out the smaller one. Finally, we filtered out deletion loci with
595 <50% average mappability (see above). This resulted in 11,390 autosomal large deletions from 661 AFR genomes.

596 **4.5.3 Deletion genotyping in ancient genomes: the raw data set**

597 We genotyped all the chosen 71 ancient genomes using the 11,390 AFR autosomal deletion data set (>1 Kbps with
598 mean 10,735 bps). We used C-score <0.5 as cutoff for calling deletions. We call this the "raw data set". In total 8,106
599 deletions were identified in at least one genome.

600 **4.6 Analyzing the raw data set of ancient deletions**

601 **4.6.1 Cluster analyses of ancient genomes**

602 We performed a principal components analysis (PCA) on the deletion copy number data set where we clustered the
603 71 ancient genomes (Figure 4A). PC1 and PC2 values were computed using the R "stats" package "prcomp" function

604 using the default parameters (R Core Team, 2020). We also generated a heatmap summarizing deletion copy numbers
605 using the R "gplots" package "heatmap.2" function (Warnes *et al.*, 2020) (Supplemental Fig. 3A). These analyses
606 revealed visible differences in deletion frequency among samples reflecting the laboratory-of-origin, and the grouping
607 of genomes into two main groups, that we called Group 1 and Group 2. We compared PC1 values among the 10
608 laboratories, and between Group 1 vs. Group 2, using the Wilcoxon rank-sum test (R "stats" package "wilcox.test"
609 function) and the Kruskal-Wallis test (R "stats" package "kruskal.test" function), respectively (R Core Team, 2020). We
610 further tested each of the 8106 deletions identified in at least one genome for differences in frequency among the 10
611 laboratories using the Kruskal-Wallis test, and corrected for multiple testing with the Benjamini-Hochberg correction
612 (R "stats" package "p.adjust" function with the "BH" parameter) (R Core Team, 2020).

613 **4.6.2 Cluster analyses of deletion events**

614 We conducted k-means clustering analysis of deletions in the "raw data set" by using the "kmeans" function offered by
615 the R "stats" package (R Core Team, 2020), and clustering deletions based on deletion copy numbers across the 71
616 genomes. We set the seed 123 for reproducibility. Analysis of within-group sum of squares did not reveal an optimal k
617 (Supplemental Fig. 6A), and therefore used both $k = 8$ and $k = 16$. After clustering, we plotted mean copy numbers of
618 the individuals in our sample set for each cluster using the R "ggplot2" package "ggplot" function (Wickham, 2016)
619 (Supplemental Fig. 4, 5).

620 We observed that Group 1 and Group 2 were different in terms of their mean copy numbers as expected. We conducted
621 the Wilcoxon rank-sum test on the mean deletion copy numbers of the two groups for each cluster (R "stats" package
622 "wilcox.test" function) (R Core Team, 2020). We also calculated effect size in mean deletion copy numbers between
623 Group 1 vs. Group 2 using a web-based effect size calculator (Becker, 2000).

624 **4.6.3 Mappability and mapping quality analysis of deletions clusters**

625 For deletion events in each k-means cluster, we compiled (i) average mappability values (see above) per deletion, and
626 (ii) mapping quality (MAPQ) values of reads mapping to those locations, separately for Group 1 and Group 2 genomes.
627 We compared these using violin plots of average mappability per deletion for each cluster ("geom_violin" function in
628 the R "ggplot2" package)(Supplemental Fig. 7, 8, 13) (Wickham, 2016). We plotted the mean difference in MAPQ
629 values between Group 1 and 2 versus the total difference in copy number between Group 1 and 2 using the "ggplot"
630 function (R "ggplot2" package) (Supplemental Fig. 14) (Wickham, 2016).

631 **4.7 Creating and analyzing the refined deletion data set and the SNP data set**

632 **4.7.1 SNP genotyping in ancient genomes**

633 Following the same reasoning as above regarding ascertainment bias, we used an African population to create a SNP
634 genotyping set for calling SNPs in the ancient genomes. Specifically, we used the 1000 Genomes Yoruba data set, which
635 included a total of 5,987,516 autosomal bi-allelic SNPs with a minor allele frequency of $\geq 10\%$ in 108 African Yoruba
636 genomes in Phase 3 of the 1000 Genomes Project (The 1000 Genomes Project Consortium, 2015). First, all reads in all
637 BAM files were clipped (trimmed) using the trimBam algorithm implemented in BamUtil (Jun *et al.*, 2015). Following
638 standard practice (Mittnik *et al.*, 2018) [REF], We trimmed (a) the end 2 bases of each read for samples prepared with
639 the Uracil-DNA-glycosylase (UDG) protocol, and (b) the end 10 bases of each read for non-UDG samples.

640 Using these BAM files of the 56 ancient individuals and the above-described list of 5,987,516 SNPs, we gener-
641 ated pseudo-haploid SNP calls at these target SNP positions by randomly selecting one read and record-
642 ing the allele carried on that read as the genotype. This was performed using the pileupCaller software
643 (<https://github.com/stschiff/sequenceTools>) on samtools mpileup output (base quality >30 and MAPQ >30) (Li *et al.*,
644 2009).

645 **4.7.2 Ancestral state determination**

646 To estimate the site frequency spectrum (SFS) for derived alleles, we polarized both deletion and SNP alleles for being
647 ancestral or derived in the human lineage. For this, we mapped loci from hg19 (GRCh37) to panTro6 (chimpanzee) and
648 to panPan2 (bonobo) using the UCSC Genome Browser tool "liftOver" with default parameters (Kent *et al.*, 2002).
649 For deletions, we filtered out deletions that did not fully map to either chimpanzee or bonobo reference genomes, as
650 these could represent derived insertions in the human lineage. The remaining deletions could thus be inferred to be
651 alleles that were derived in humans. For SNPs, we removed the positions not represent in either chimpanzee or bonobo
652 reference genomes and assigned the ancestral state as the Pan allele, only if both chimpanzee and bonobo carried same
653 allele. This left us with 4,907,535 SNP positions with derived allele information.

654 **4.7.3 Creating the refined deletion data set**

655 We filtered deletions in the raw data set for high mappability (≥ 0.9 average mappability), which left us with 10,018
656 deletions. We further removed 15 genomes identified as outliers in the k-means clustering analyses. We then selected
657 reads with >30 mapping quality (MAPQ) in all BAM files, and we genotyped the 10,018 AFR deletions in the
658 remaining 56 genomes. We call this the "refined data set". In addition, we created a subset of this data set than only
659 includes the 9,579 deletions derived in the human lineage (see above). After refining our data set, we also checked
660 its general properties. We plotted size distribution in logarithmic scale, deletion allele frequency distribution and
661 relative frequency distribution among observed heterozygous deletions over homozygous deletions using R's "graphics"
662 package hist function (Supplemental Fig. 15) (R Core Team, 2020). We also plotted relative deletion (homozygous or
663 heterozygous) frequencies of 9579 deletions for each individual in our refined data set using R's "graphics" package
664 matplot function (R Core Team, 2020).

665 **4.7.4 Site frequency spectrum calculation for deletions and SNPs**

666 Here our goal was to compare the SFS across deletions and SNPs called in ancient genomes. Because the ancient SNP
667 genotypes are pseudo-haploidized, we performed the same pseudo-haploidization process on the deletion data set. For
668 this, for any heterozygous call in the deletion data set, we randomly assigned either of the homozygous states, using the
669 R "sample" function (i.e., we converted 1's to 0's or 2's with 50% chance). We then counted derived alleles at each
670 locus, for deletions and for SNPs, and divided by the total number of genomes where an allele was observed at that
671 locus (i.e., removing the missing data). We plotted the site-frequency spectrum analysis on both deletions and SNPs
672 using R's "ggplot2" package geom_histogram function (Wickham, 2016). We also calculated the correlation between
673 the deletion size in logarithmic scale and the frequency using R's "stats" package cor.test with method parameter set to
674 "s" indicating Spearman's correlation (R Core Team, 2020).

675 **4.7.5 Evolutionary conservation**

676 To measure evolutionary conservation for genes that overlapped deletions, we retrieved non-synonymous (dN) and
677 synonymous (dS) substitution rate estimates between human (GRCh37) and the mouse genome (GRCm38) per gene
678 from Ensembl (v75) using the BiomaRt tool. We queried 18,112 genes with dN, dS values and calculated the dN/dS
679 ratio (or Ka/Ks) per gene. The ratio for genes with more than one dN or dS values were calculated as the mean dN or
680 dS per gene. We then intersected our deletions with the genes with dN/dS values using BEDTools (Quinlan and Hall,
681 2010), which yielded $n = 401$ genes. We divided the deletions in our data set into two groups by the deletion allele
682 frequency per gene: high versus low frequency. We calculated the mean frequency among the 401 genes and set this
683 value as a threshold for high and low frequency. The frequencies above the threshold were considered high and vice
684 versa. We plotted the dN/dS ratios of the groups defined above using the R package "ggplot2" and the "geom_boxplot"
685 function (Wickham, 2016).

686 **4.8 Genetic distance analyses using deletions and SNPs**

687 Here our goal was to calculate overall genetic distances among the 56 ancient genomes using deletion allele frequencies
688 and using SNPs, and further to compare the distances. We calculated distances using the commonly used outgroup-f3
689 statistics, which measures shared genetic drift between two samples relative to an outgroup, and is implemented as
690 qp3pop in Admixtools v.7.0 (Patterson *et al.*, 2012). The outgroup-f3 values were calculated for each pair of 56
691 individuals (a) in the deletion and (b) in the SNP data sets, using the African Yoruba as outgroup in both cases. To
692 convert the deletion data set to eigenstrat format, which Admixtools requires, we encoded the first nucleotide of each
693 deletion as the reference allele, and the alternative allele was randomly assigned among the remaining 3 nucleotides
694 using custom python script. We thus calculated a pairwise similarity matrix for both data sets. Genetic distances
695 were calculated as $1-f_3$. Distances were then summarized using multidimensional scaling (MDS) with the "cmdscale"
696 function of R (R Core Team, 2020).

697 We further performed the Mantel test to compare the f3-based similarity matrices calculated using SNPs and deletions.
698 We used the "mantel" function in the R-package "vegan" with parameter "method=spearman" (Oksanen *et al.*, 2013).

699 **Data access**

700 CONGA is implemented in C programming language and its source code is available under BSD 3-clause license at
701 <https://github.com/asylvz/CONGA>. Simulated datasets and predictions of each tool can be accessed through Zenodo
702 (10.5281/zenodo.5555990)

703 **Competing interest statement**

704 The authors declare no competing interests.

705 **Acknowledgements**

706 The authors would like to thank Gözde Zeliha Turan for her suggestions with mappability data and Kivılcım Başak
707 Vural for technical support. This work was supported by the ERC Consolidator grant “NEOGENE” (Project No.:
708 772390).

709 **Author Contributions**

710 AS developed and implemented the algorithm, performed simulations and down-sampling experiments. SSÇ and DK
711 conducted technical and evolutionary analyses on real data. CA contributed to algorithm design. MS led the project and
712 coordinated the activities. All authors contributed to editing the manuscript and participated in weekly discussions.

713 **References**

- 714 (2019). Picard toolkit. <https://broadinstitute.github.io/picard/>.
- 715 Abyzov, A. *et al.* (2011). CNVnator: an approach to discover, genotype, and characterize typical and atypical CNVs
716 from family and population genome sequencing. *Genome Res*, **21**(6), 974–984.
- 717 Alkan, C. (2020). Automatic characterization of copy number polymorphism using high throughput sequencing. *Turkish*
718 *Journal of Electrical Engineering & Computer Sciences*, **28**(1), 253–261.
- 719 Alkan, C. *et al.* (2009). Personalized copy number and segmental duplication maps using next-generation sequencing.
720 *Nat Genet*, **41**(10), 1061–1067.
- 721 Alkan, C. *et al.* (2011). Genome structural variation discovery and genotyping. *Nat Rev Genet*, **12**(5), 363–376.
- 722 Allentoft, M. E. *et al.* (2015). Population genomics of bronze age eurasia. *Nature*, **522**(7555), 167.
- 723 Almarri, M. A. *et al.* (2020). Population structure, stratification, and introgression of human structural variation. *Cell*,
724 **182**(1), 189–199.
- 725 Antonio, M. L. *et al.* (2019). Ancient rome: A genetic crossroads of europe and the mediterranean. *Science*, **366**(6466),
726 708–714.
- 727 Audano, P. A. *et al.* (2019). Characterizing the major structural variant alleles of the human genome. *Cell*, **176**(3),
728 663–675.
- 729 Becker, L. A. (2000). Effect size (es).
- 730 Bergström, A. *et al.* (2020). Origins and genetic legacy of prehistoric dogs. *Science*, **370**(6516), 557–564.
- 731 Bhattacharya, S. *et al.* (2018). Whole-genome sequencing of atacama skeleton shows novel mutations linked with
732 dysplasia. *Genome research*, **28**(4), 423–431.
- 733 Boeva, V. *et al.* (2012). Control-FREEC: a tool for assessing copy number and allelic content using next-generation
734 sequencing data. *Bioinformatics (Oxford, England)*, **28**, 423–425.
- 735 Broushaki, F. *et al.* (2016). Early neolithic genomes from the eastern fertile crescent. *Science*, **353**(6298), 499–503.
- 736 Chaisson, M. J. P. *et al.* (2015). Resolving the complexity of the human genome using single-molecule sequencing.
737 *Nature*, **517**, 608–611.
- 738 Chaisson, M. J. P. *et al.* (2019). Multi-platform discovery of haplotype-resolved structural variation in human genomes.
739 *Nature Communications*, **10**, 1784.
- 740 Chan, Y. F. *et al.* (2010). Adaptive evolution of pelvic reduction in sticklebacks by recurrent deletion of a *pitx1* enhancer.
741 *science*, **327**(5963), 302–305.

- 742 Chen, X. *et al.* (2016). Manta: rapid detection of structural variants and indels for germline and cancer sequencing
743 applications. *Bioinformatics*, **32**, 1220–1222.
- 744 Clark, A. G. *et al.* (2005). Ascertainment bias in studies of human genome-wide polymorphism. *Genome research*,
745 **15**(11), 1496–1502.
- 746 Collins, R. L. *et al.* (2020). A structural variation reference for medical and population genetics. *Nature*, **581**(7809),
747 444–451.
- 748 Conrad, D. F. and Hurler, M. E. (2007). The population genetics of structural variation. *Nature genetics*, **39**(7),
749 S30–S36.
- 750 Conrad, D. F. *et al.* (2010). Origins and functional impact of copy number variation in the human genome. *Nature*,
751 **464**(7289), 704–712.
- 752 de Barros Damgaard, P. *et al.* (2018a). 137 ancient human genomes from across the eurasian steppes. *Nature*, **557**(7705),
753 369–374.
- 754 de Barros Damgaard, P. *et al.* (2018b). The first horse herders and the impact of early bronze age steppe expansions
755 into asia. *Science*, **360**(6396), eaar7711.
- 756 Derrien, T. *et al.* (2012). Fast computation and applications of genome mappability. *PloS one*, **7**(1), e30377.
- 757 Eisfeldt, J. *et al.* (2017). TIDDIT, an efficient and comprehensive structural variant caller for massive parallel sequencing
758 data. *F1000Research*, **6**, 664.
- 759 ENCODE Project Consortium (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature*,
760 **489**, 57–74.
- 761 Frantz, L. A. *et al.* (2020). Animal domestication in the era of ancient genomics. *Nature Reviews Genetics*, pages 1–12.
- 762 Fu, Q. *et al.* (2014). Genome sequence of a 45,000-year-old modern human from western siberia. *Nature*, **514**(7523),
763 445–449.
- 764 Gamba, C. *et al.* (2014). Genome flux and stasis in a five millennium transect of european prehistory. *Nature*
765 *communications*, **5**(1), 1–9.
- 766 Girirajan, S. *et al.* (2011). Human copy number variation and complex genetic disease. *Annual review of genetics*, **45**,
767 203–226.
- 768 Gonzalez, E. *et al.* (2005). The influence of CCL3L1 gene-containing segmental duplications on HIV-1/AIDS
769 susceptibility. *Science*, **307**(5714), 1434–1440.
- 770 González-Fortes, G. *et al.* (2017). Paleogenomic evidence for multi-generational mixing between neolithic farmers and
771 mesolithic hunter-gatherers in the lower danube basin. *Current Biology*, **27**(12), 1801–1810.
- 772 Günther, T. *et al.* (2015). Ancient genomes link early farmers from atapuerca in spain to modern-day basques.
773 *Proceedings of the National Academy of Sciences*, **112**(38), 11917–11922.
- 774 Haber, M. *et al.* (2017). Continuity and admixture in the last five millennia of levantine history from ancient canaanite
775 and present-day lebanese genome sequences. *The American Journal of Human Genetics*, **101**(2), 274–282.
- 776 Haber, M. *et al.* (2019). A transient pulse of genetic admixture from the crusaders in the near east identified from
777 ancient genome sequences. *The American Journal of Human Genetics*, **104**(5), 977–984.
- 778 Hardwick, R. J. *et al.* (2011). A worldwide analysis of beta-defensin copy number variation suggests recent selection of
779 a high-expressing defb103 gene copy in east asia. *Human mutation*, **32**(7), 743–750.
- 780 Ho, S. S. *et al.* (2020). Structural variation in the sequencing era. *Nature Reviews Genetics*, **21**(3), 171–189.
- 781 Hofmanová, Z. *et al.* (2016). Early farmers from across europe directly descended from neolithic aegeans. *Proceedings*
782 *of the National Academy of Sciences*, **113**(25), 6886–6891.
- 783 Hsieh, P. *et al.* (2019). Adaptive archaic introgression of copy number variants and the discovery of previously unknown
784 human genes. *Science*, **366**(6463).

- 785 Jones, E. R. *et al.* (2015). Upper palaeolithic genomes reveal deep roots of modern eurasians. *Nature communications*,
786 **6**(1), 1–8.
- 787 Jones, E. R. *et al.* (2017). The neolithic transition in the baltic was not driven by admixture with early european farmers.
788 *Current Biology*, **27**(4), 576–582.
- 789 Jun, G. *et al.* (2015). An efficient and scalable analysis framework for variant extraction and refinement from
790 population-scale dna sequence data. *Genome research*, **25**(6), 918–925.
- 791 Karakoc, E. *et al.* (2012). Detection of structural variants and indels within exome data. *Nat Methods*, **9**(2), 176–178.
- 792 Karimzadeh, M. *et al.* (2018). Umap and Bimap: quantifying genome and methylome mappability. *Nucleic Acids*
793 *Research*, **46**(20), e120–e120.
- 794 Keller, A. *et al.* (2012). New insights into the tyrolean iceman’s origin and phenotype as inferred by whole-genome
795 sequencing. *Nature communications*, **3**(1), 1–9.
- 796 Kent, W. J. *et al.* (2002). The human genome browser at ucsc. *Genome research*, **12**(6), 996–1006.
- 797 Kılınc, G. M. *et al.* (2016). The demographic development of the first farmers in anatolia. *Current Biology*, **26**(19),
798 2659–2666.
- 799 Kircher, M. (2012). Analysis of high-throughput ancient dna sequencing data. In *Ancient DNA*, pages 197–228.
800 Springer.
- 801 Koehler, R. *et al.* (2011). The uniqueome: a mappability resource for short-tag sequencing. *Bioinformatics*, **27**(2),
802 272–274.
- 803 Kothapalli, K. S. *et al.* (2016). Positive selection on a regulatory insertion–deletion polymorphism in fads2 influences
804 apparent endogenous synthesis of arachidonic acid. *Molecular biology and evolution*, **33**(7), 1726–1739.
- 805 Kousathanas, A. *et al.* (2017). Inferring heterozygosity from ancient and low coverage genomes. *Genetics*, **205**(1),
806 317–332.
- 807 Krzewińska, M. *et al.* (2018). Ancient genomes suggest the eastern pontic-caspian steppe as the source of western iron
808 age nomads. *Science advances*, **4**(10), eaat4457.
- 809 Layer, R. M. *et al.* (2014). LUMPY: a probabilistic framework for structural variant discovery. *Genome Biol*, **15**(6),
810 R84.
- 811 Lazaridis, I. *et al.* (2014). Ancient human genomes suggest three ancestral populations for present-day europeans.
812 *Nature*, **513**(7518), 409–413.
- 813 Levy-Sakin, M. *et al.* (2019). Genome maps across 26 human populations reveal population-specific patterns of
814 structural variation. *Nature communications*, **10**, 1025.
- 815 Li, H. and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*,
816 **25**(14), 1754–1760.
- 817 Li, H. *et al.* (2009). The sequence alignment/map format and SAMtools. *Bioinformatics*, **25**(16), 2078–2079.
- 818 Link, V. *et al.* (2017). Atlas: analysis tools for low-depth and ancient samples. *BioRxiv*, page 105346.
- 819 Llorente, M. G. *et al.* (2015). Ancient ethiopian genome reveals extensive eurasian admixture in eastern africa. *Science*,
820 **350**(6262), 820–822.
- 821 Marciniak, S. and Perry, G. H. (2017). Harnessing ancient genomes to study the history of human adaptation. *Nature*
822 *Reviews Genetics*, **18**(11), 659–674.
- 823 Martiniano, R. *et al.* (2017). The population genomics of archaeological transition in west iberia: Investigation of
824 ancient substructure using imputation and haplotype-based methods. *PLoS genetics*, **13**(7), e1006852.
- 825 Mathieson, S. and Mathieson, I. (2018). Fads1 and the timing of human adaptation to agriculture. *Molecular biology*
826 *and evolution*, **35**(12), 2957–2970.

- 827 McLean, C. Y. *et al.* (2011). Human-specific loss of regulatory dna and the evolution of human-specific traits. *Nature*,
828 **471**(7337), 216–219.
- 829 Miller, C. A. *et al.* (2011). ReadDepth: a parallel R package for detecting copy number alterations from short sequencing
830 reads. *PloS one*, **6**, e16327.
- 831 Mittnik, A. *et al.* (2018). The genetic prehistory of the baltic sea region. *Nature communications*, **9**(1), 1–11.
- 832 Mu, J. C. *et al.* (2015). VarSim: a high-fidelity simulation and validation framework for high-throughput genome
833 sequencing with cancer applications. *Bioinformatics*, **31**(9), 1469–1471.
- 834 Nuttle, X. *et al.* (2016). Emergence of a homo sapiens-specific gene family and chromosome 16p11.2 cnv susceptibility.
835 *Nature*, **536**(7615), 205–209.
- 836 Oksanen, J. *et al.* (2013). Package ‘vegan’. *Community ecology package, version*, **2**(9), 1–295.
- 837 Olalde, I. *et al.* (2014). Derived immune and ancestral pigmentation alleles in a 7,000-year-old mesolithic european.
838 *Nature*, **507**(7491), 225–228.
- 839 Oliva, A. *et al.* (2021). Bwa-mem is not the best aligner for ancient dna short reads. *bioRxiv*.
- 840 Patterson, N. *et al.* (2012). Ancient admixture in human history. *Genetics*, **192**(3), 1065–1093.
- 841 Pedersen, J. S. *et al.* (2014). Genome-wide nucleosome map and cytosine methylation levels of an ancient human
842 genome. *Genome research*, **24**(3), 454–466.
- 843 Perry, G. H. *et al.* (2007). Diet and the evolution of human amylase gene copy number variation. *Nature genetics*,
844 **39**(10), 1256–1260.
- 845 Pockrandt, C. *et al.* (2020). Genmap: ultra-fast computation of genome mappability. *Bioinformatics*.
- 846 Prüfer, K. (2018). snpad: An ancient dna genotype caller. *Bioinformatics*, **34**(24), 4165–4171.
- 847 Quinlan, A. R. and Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features.
848 *Bioinformatics*, **26**(6), 841–842.
- 849 R Core Team (2020). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical
850 Computing, Vienna, Austria.
- 851 Raghavan, M. *et al.* (2014). Upper palaeolithic siberian genome reveals dual ancestry of native americans. *Nature*,
852 **505**(7481), 87–91.
- 853 Rasmussen, M. *et al.* (2010). Ancient human genome sequence of an extinct palaeo-eskimo. *Nature*, **463**(7282),
854 757–762.
- 855 Rasmussen, M. *et al.* (2014). The genome of a late pleistocene human from a clovis burial site in western montana.
856 *Nature*, **506**(7487), 225–229.
- 857 Rausch, T. *et al.* (2012). DELLY: structural variant discovery by integrated paired-end and split-read analysis.
858 *Bioinformatics*, **28**(18), i333–i339.
- 859 Renaud, G. *et al.* (2017). gargammel: a sequence simulator for ancient dna. *Bioinformatics*, **33**(4), 577–579.
- 860 Rogers, R. L. and Slatkin, M. (2017). Excess of genomic defects in a woolly mammoth on wrangel island. *PLoS*
861 *genetics*, **13**(3), e1006601.
- 862 Rohland, N. *et al.* (2015). Partial uracil–dna–glycosylase treatment for screening of ancient dna. *Philosophical*
863 *Transactions of the Royal Society B: Biological Sciences*, **370**(1660), 20130624.
- 864 Saitou, M. and Gokcumen, O. (2020). An evolutionary perspective on the impact of genomic copy number variation on
865 human health. *Journal of molecular evolution*, **88**(1), 104–119.
- 866 Schubert, M. *et al.* (2016). Adapterremoval v2: rapid adapter trimming, identification, and read merging. *BMC research*
867 *notes*, **9**(1), 1–7.
- 868 Sedlazeck, F. J. *et al.* (2018). Accurate detection of complex structural variations using single-molecule sequencing.
869 *Nature methods*, **15**, 461–468.

- 870 Seguin-Orlando, A. *et al.* (2014). Genomic structure in europeans dating back at least 36,200 years. *Science*, **346**(6213),
871 1113–1118.
- 872 Shapiro, á. and Hofreiter, M. (2014). A paleogenomic perspective on evolution and gene function: new insights from
873 ancient dna. *Science*, **343**(6169), 1236573.
- 874 Sikora, M. *et al.* (2019). The population history of northeastern siberia since the pleistocene. *Nature*, **570**(7760),
875 182–188.
- 876 Skoglund, P. and Mathieson, I. (2018). Ancient genomics of modern humans: the first decade. *Annual review of*
877 *genomics and human genetics*, **19**, 381–404.
- 878 Smith, D. R. *et al.* (2008). Rapid whole-genome mutational profiling using next-generation sequencing technologies.
879 *Genome Res*, **18**(10), 1638–1642.
- 880 Smith, S. D. *et al.* (2015). Grom-rd: resolving genomic biases to improve read depth detection of copy number variants.
881 *PeerJ*, **3**, e836.
- 882 Smith, S. D. *et al.* (2017). Evolutionary adaptation revealed by comparative genome analysis of woolly mammoths and
883 elephants. *DNA Research*, **24**(4), 359–369.
- 884 Soylev, A. *et al.* (2017). Toolkit for automated and rapid discovery of structural variants. *Methods*, **129**, 3–7.
- 885 Soylev, A. *et al.* (2019). Discovery of tandem and interspersed segmental duplications using high-throughput sequencing.
886 *Bioinformatics*, **35**, 3923–3930.
- 887 Stankiewicz, P. and Lupski, J. R. (2010). Structural variation in the human genome and its role in disease. *Annu Rev*
888 *Med*, **61**, 437–455.
- 889 Sudmant, P. H. *et al.* (2010). Diversity of human copy number variation and multicopy genes. *Science*, **330**(6004),
890 641–646.
- 891 Sudmant, P. H. *et al.* (2015). An integrated map of structural variation in 2,504 human genomes. *Nature*, **526**(7571),
892 75–81.
- 893 The 1000 Genomes Project Consortium (2015). A global reference for human genetic variation. *Nature*, **526**(7571),
894 68–74.
- 895 Vieira, F. G. *et al.* (2016). Estimating ibd tracts from low coverage ngs data. *Bioinformatics*, **32**(14), 2096–2102.
- 896 Warnes, G. R. *et al.* (2020). *gplots: Various R Programming Tools for Plotting Data*. R package version 3.1.1.
- 897 Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.
- 898 Xue, Y. *et al.* (2008). Adaptive evolution of ugt2b17 copy-number variation. *The American Journal of Human Genetics*,
899 **83**(3), 337–346.
- 900 Yaka, R. *et al.* (2021). Variable kinship patterns in neolithic anatolia revealed by ancient genomes. *Current Biology*,
901 **31**(11), 2455–2468.
- 902 Zhang, F. *et al.* (2009). Copy number variation in human health, disease, and evolution. *Annual review of genomics and*
903 *human genetics*, **10**, 451–481.
- 904 Zook, J. M. *et al.* (2020). A robust benchmark for detection of germline large deletions and insertions. *Nature*
905 *Biotechnology*, pages 1–9.