

Supplementary Materials:
BADGER: evaluating the performance of ancient DNA genetic
relatedness estimation methods using high-fidelity pedigree simulations.

001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046

047	Contents	
048		
049		
050	1 Supplementary Figures	4
051		
052	1.1 Figure S1: Alternate template pedigree including inbred individuals . .	4
053		
054	1.2 Figure S2: Alternative input pedigrees used to evaluate the impact	
055	of inbreeding on the normalisation procedure of correctKin, KIN and	
056	READv2.	5
057		
058	1.3 Figure S3: Comparing the performance of READv2 against its prede-	
059	cessor, READ	6
060		
061	1.4 Figure S4: Accuracy and bias of r-coefficients as a function of sequencing	
062	depth (pmd-mask).	7
063		
064	1.5 Figure S5: Accuracy and bias of r-coefficients as a function of sequencing	
065	depth (mapDamage2)	8
066		
067	1.6 Figure S6: Full grid of confusion matrices and UOC values across	
068	increasing values of sequencing depth (mapDamage2)	9
069		
070	1.7 Figure S7: Accuracy and bias of r-coefficients as a function of contam-	
071	ination rate (AFR)	10
072		
073	1.8 Figure S8: Accuracy and bias of r-coefficients as a function of contam-	
074	ination (GBR)	11
075		
076	1.9 Figure S9: Full grid of confusion matrices and UOC values across	
077	increasing values of contamination (AFR)	12
078		
079	1.10 Figure S10: Full grid of confusion matrices and UOC values across	
080	increasing values of contamination (GBR)	13
081		
082	1.11 Figure S11: Impact of admixture on the accuracy and bias of r-coefficients	14
083		
084	1.12 Figure S12: Confusion matrices and UOC values across increasing values	
085	of sequencing depth (ASW)	15
086		
087	1.13 Figure S13: Ancestry proportions of admixed American populations. .	16
088		
089		
090		
091		
092		

1.14	Figure S14: Average heterozygosity rate of the European CEU population, and admixed American populations.	17	093 094 095
1.15	Figure S15: Accuracy and bias of r-coefficients estimates for pairwise comparisons involving inbred individuals	18	096 097 098 099
1.16	Figure S16: Impact of inbreeding on the accuracy and bias of r-coefficients estimates for pairwise comparisons involving outbred individuals	19	100 101 102 103 104
2	Material and Methods	20	105 106
2.1	Description of BADGER’s simulation pipeline	20	107 108
2.1.1	1000 genomes dataset pre-processing	21	109
2.1.2	Pedigree simulations	21	110 111
2.1.3	Ancient DNA simulations	22	112 113
2.1.4	Alignment	23	114
2.1.5	Quality filtering and preprocessing of alignment files	24	115 116
2.1.6	Correction of <i>post-mortem</i> deaminations	24	117 118
2.1.7	Variant calling	24	119 120
2.1.8	Genetic relatedness estimation	25	121
2.2	Statistical analysis and benchmark using badger.plots	28	122 123
2.2.1	Estimation of classification performance	29	124
2.2.2	Average accuracy and bias of relatedness coefficients	30	125 126
2.3	Key Resources Table	32	127 128 129
3	Description of the pmd-mask command line utility	34	130
3.1	Rationale, behaviour and workflow description	34	131 132
3.2	Pseudo-code describing the main algorithm of pmd-mask	37	133 134 135 136 137 138

1 Supplementary Figures

1.1 Figure S1: Alternate template pedigree including inbred individuals

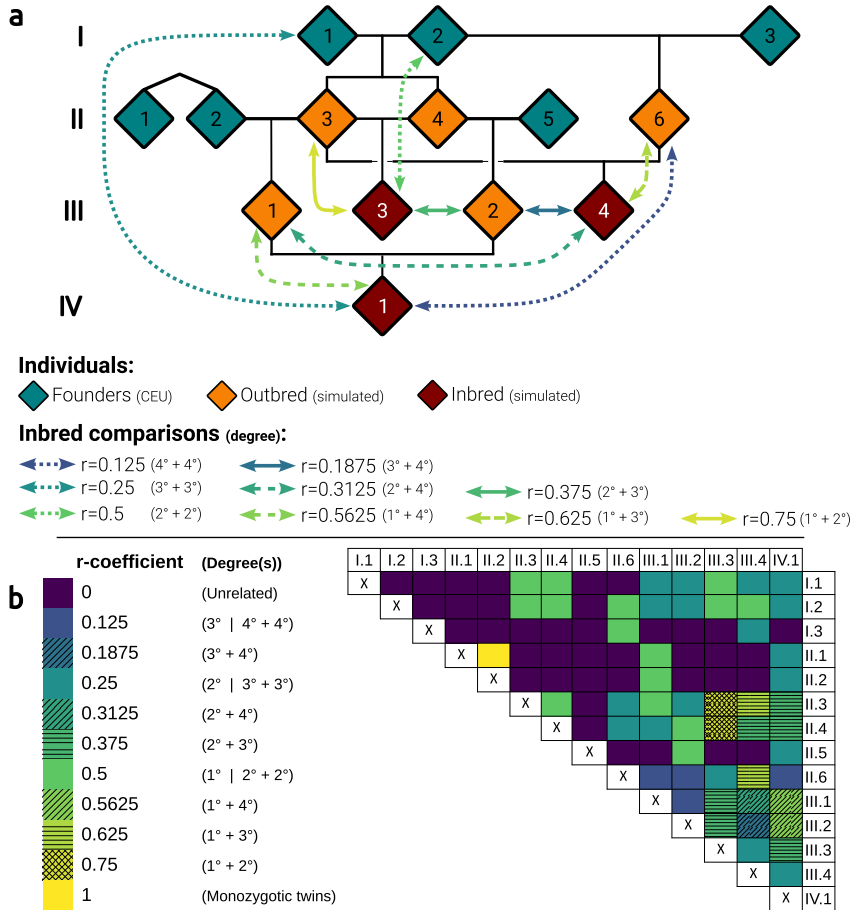


Fig. S1: Alternate template pedigree including three inbred individuals. **a:** Diagram of the alternative template pedigree, provided to BADGER when evaluating the impact of close inbreeding. Coloured arrows denote the pairwise relationships investigated during the inbreeding benchmark. Additional inbred individuals are coloured in dark red, with individuals III.3, III.4 and IV.1 being the result of a mating respectively involving two siblings, two half-siblings, and two first-cousins. **b:** Pairwise matrix of the relationships defined by the simulated pedigree shown in (a). A complete description of the pairwise relationships contained within the input pedigrees used throughout this study is described in Supplementary Table S2

1.2 Figure S2: Alternative input pedigrees used to evaluate
the impact of inbreeding on the normalisation procedure
of correctKin, KIN and READv2.

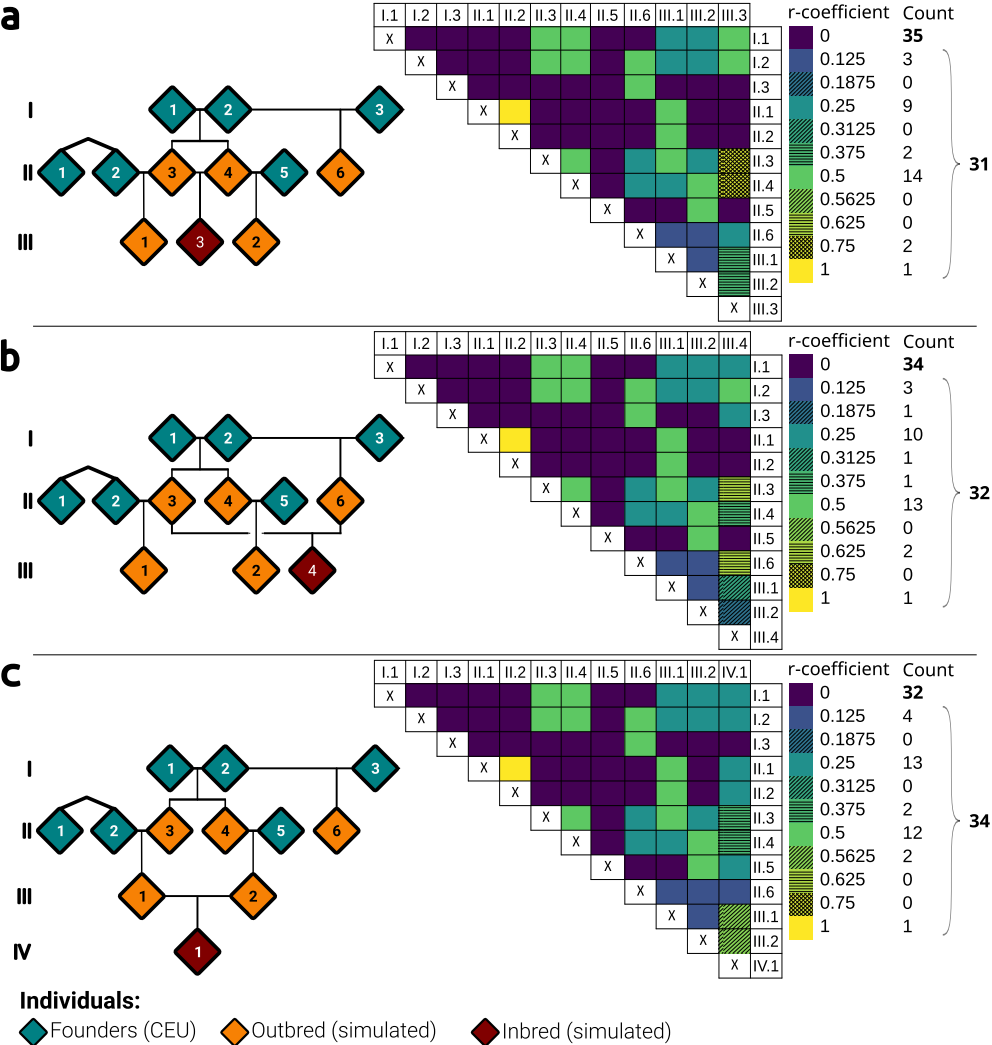


Fig. S2: Diagram and pairwise relationship matrix of three alternative pedigree topologies, given as input to correctKin, KIN and READv2 methods, where only one out of the three inbred individuals (III.3, III.4 and IV.1) is included in the tested cohort. **a:** Full-siblings scenario. **b:** Half-siblings scenario. **c:** First-cousins scenario.

1.3 Figure S3: Comparing the performance of READv2 against its predecessor, READ

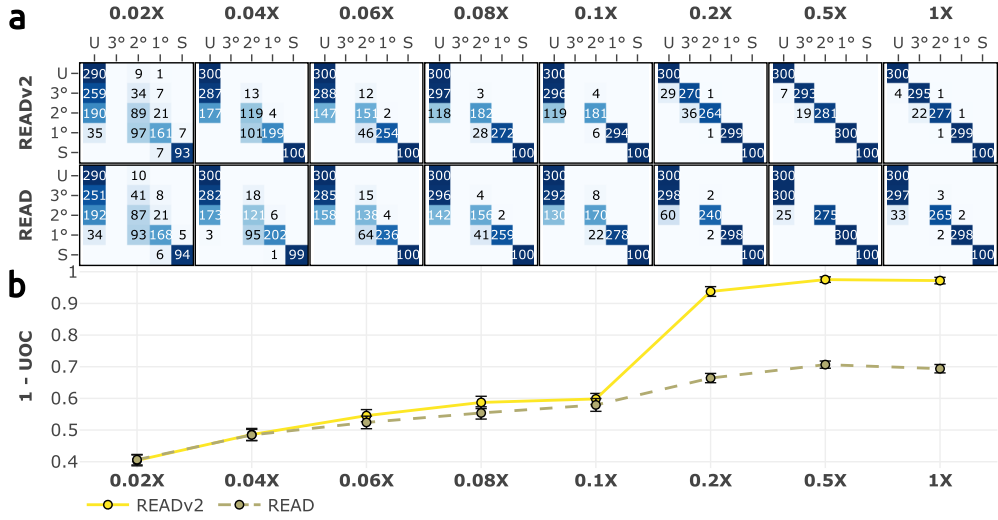


Fig. S3: Benchmark results across increasing values of sequencing depth for the READ (dashed tan line) method and its updated version, READv2 (solid yellow line). **a:** Confusion matrices of the READ and READv2 methods, confronting expected and predicted relationships. Expected and predicted values are displayed in rows and columns, respectively. 1°, 2°, 3° correspond to first-, second-, and third-degree relationships, respectively, *U* corresponds to "unrelated individuals", and *S* to "self" (monozygotic twins). **b:** UOC values summarizing the classification performance of each method for the considered sequencing depths. Higher values of 1 - UOC indicate higher performance.

1.4 Figure S4: Accuracy and bias of r-coefficients as a function
of sequencing depth (pmd-mask).

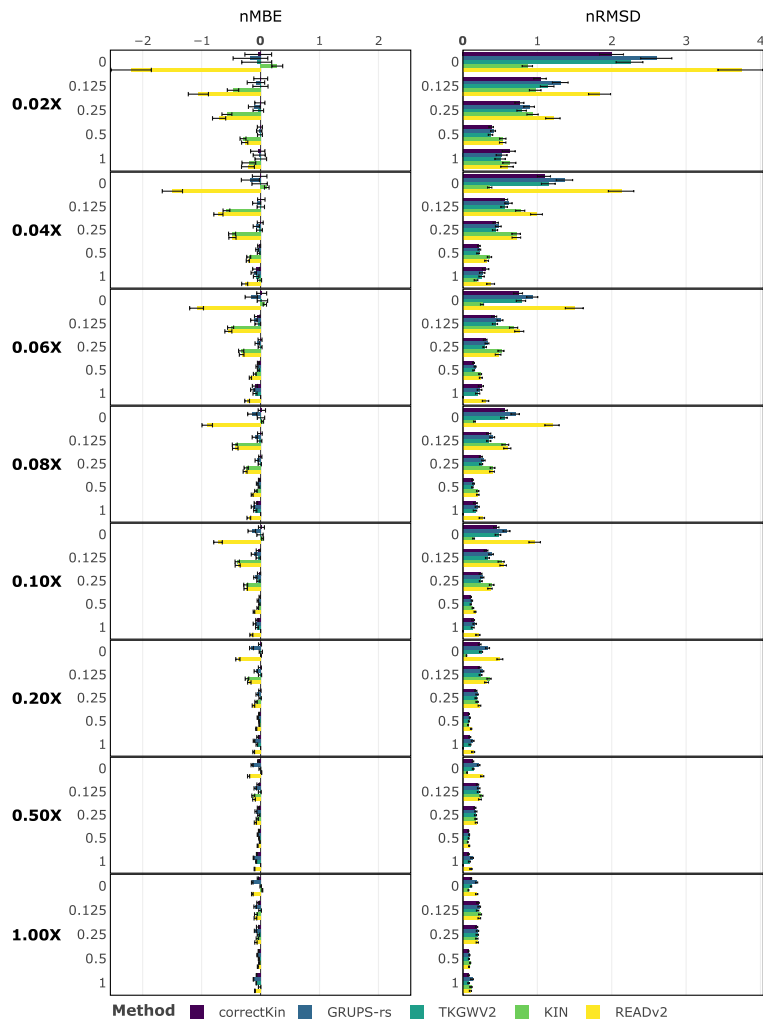


Fig. S4: Normalized estimates of MBE (left column) and $RMSD$ (right column), across all evaluated methods (bar colours), sequencing depths (rows), and expected relatedness coefficients (y-axis ticks), using sample alignment files processed through `pmd-mask`. Increasing values of $nRMSD$ indicate lower accuracy when estimating relatedness coefficients. $nMBE$ values that deviate furthest from zero indicate higher bias, with positive and negative values highlighting a tendency towards over- or under-estimating r-coefficients, respectively. Error bars represent $CI_{95\%}$ for the given estimate.

1.5 Figure S5: Accuracy and bias of r-coefficients as a function of sequencing depth (mapDamage2)

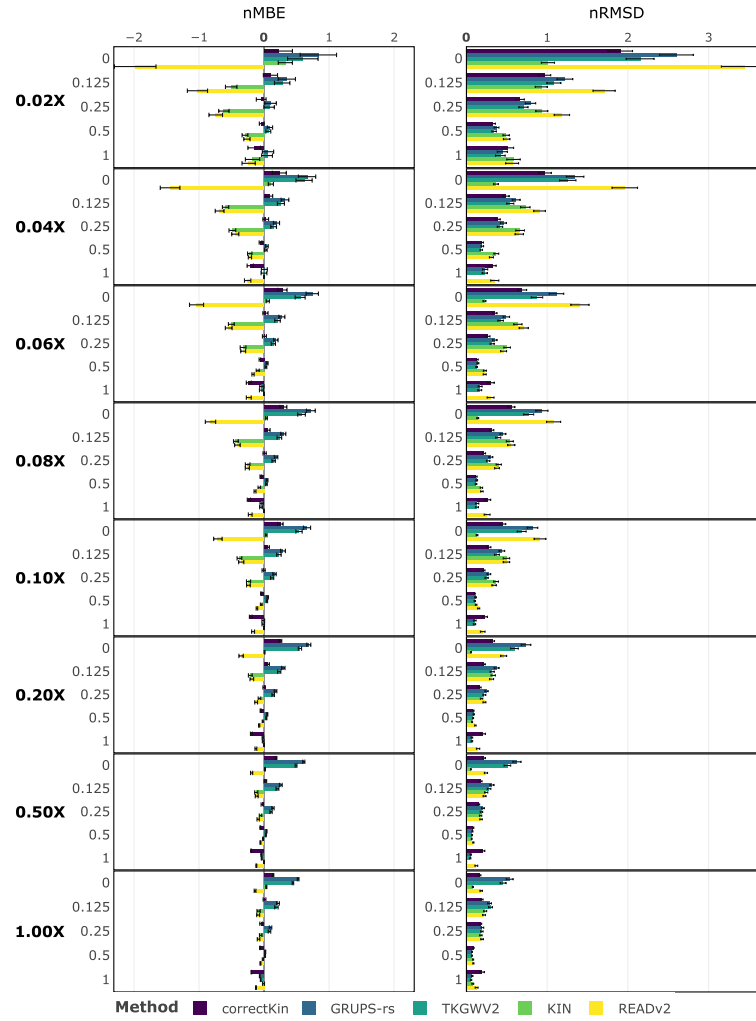


Fig. S5: *nMBE* (left column) and *nRMSD* (right column) estimates, across all evaluated methods (bar colours), sequencing depths (rows), and expected relatedness coefficients (y-axis ticks), using sample alignment files processed through `mapDamage2`. Increasing values of *nRMSD* indicate lower accuracy when estimating relatedness coefficients. *nMBE* values that deviate furthest from zero indicate higher bias, with positive and negative values highlighting a tendency towards over- or under-estimating r-coefficients, respectively. Error bars represent $CI_{95\%}$ for the given estimate.

1.6 Figure S6: Full grid of confusion matrices and UOC values
across increasing values of sequencing depth
(mapDamage2)

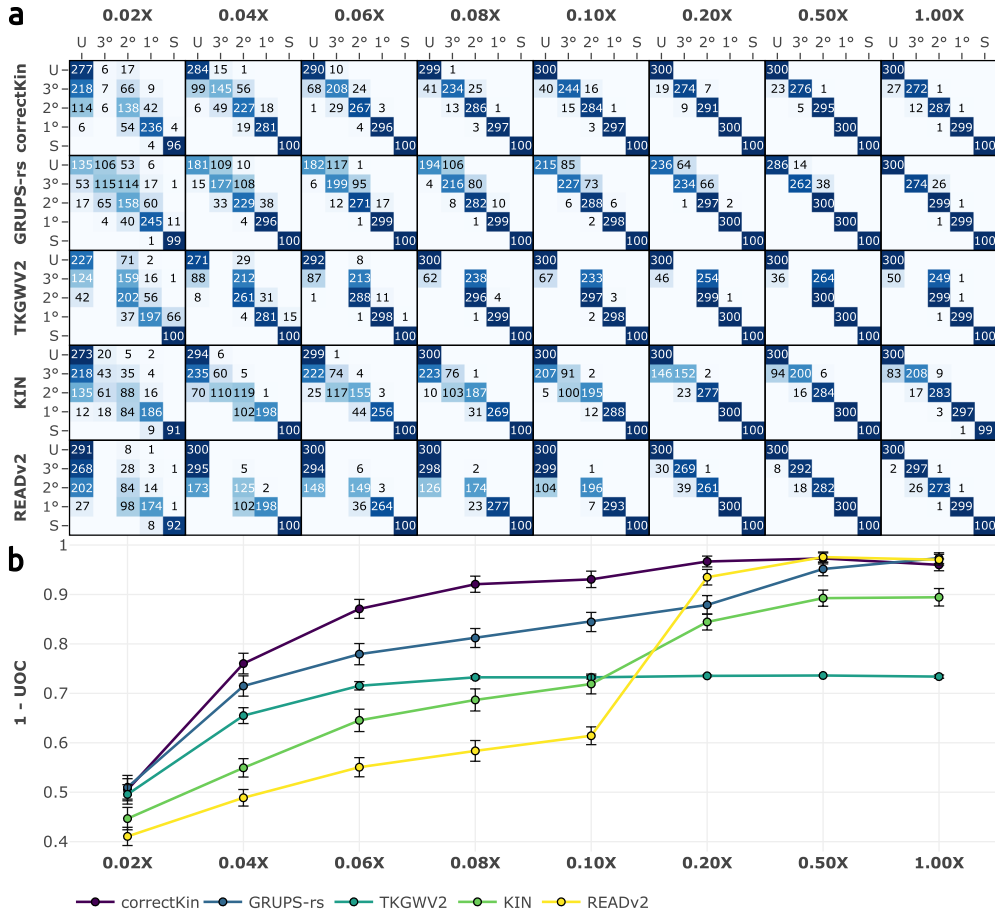


Fig. S6: Benchmark results across increasing values of sequencing depth, using sample alignment files processed through mapDamage2 *post-mortem* damage rescaling software. **a:** Confusion matrices of the five tested methods confronting expected and predicted relationships. Expected and predicted values are displayed in rows and columns, respectively. 1°, 2°, 3° correspond to first-, second-, and third-degree relationships, respectively, *U* corresponds to "unrelated individuals", and *S* to "self" (monozygotic twins). **b:** UOC values summarizing the classification performance of each method for the considered sequencing depths. Higher values of $1 - UOC$ indicate higher performance.

1.7 Figure S7: Accuracy and bias of r-coefficients as a function of contamination rate (AFR)

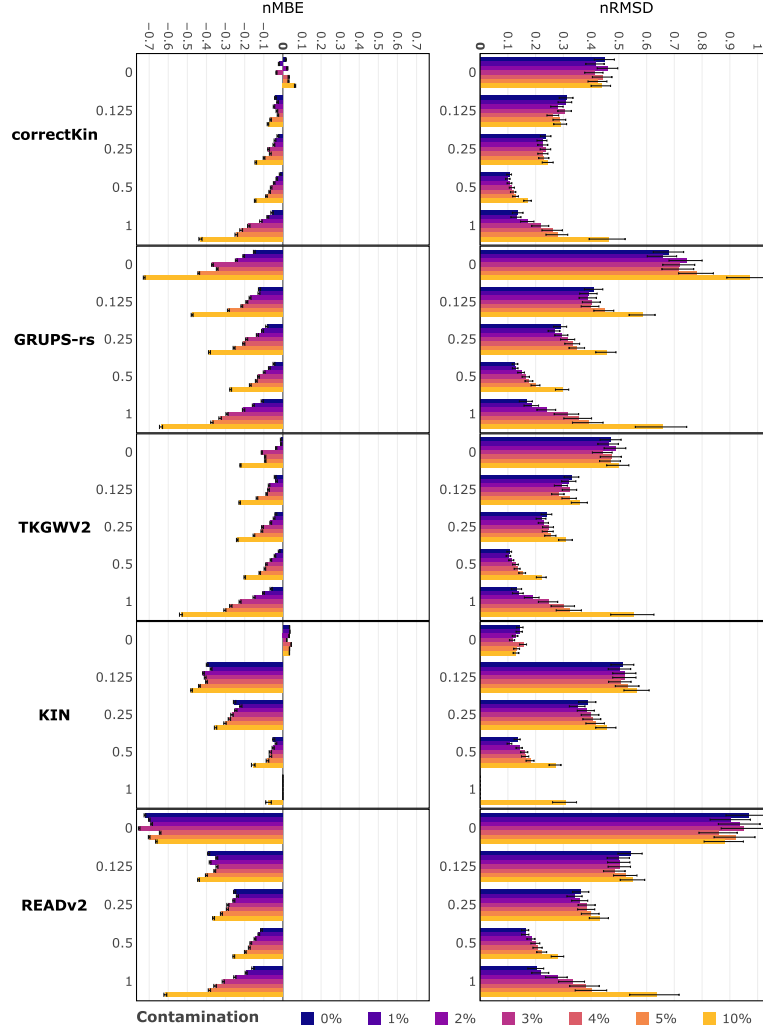


Fig. S7: *nMBE* (left column) and *nRMSD* (right column) estimates, across all evaluated methods (rows), rates of contamination from an *AFR* individual (bar colours [1], and expected relatedness coefficients (y-axis ticks). Increasing values of *nRMSD* indicate lower accuracy when estimating relatedness coefficients. *nMBE* values that deviate furthest from zero indicate higher bias, with positive and negative values highlighting a tendency towards over- or under-estimating r-coefficients, respectively. Error bars represent $CI_{95\%}$ for the given estimate.

1.8 Figure S8: Accuracy and bias of r-coefficients as a function
of contamination (GBR)

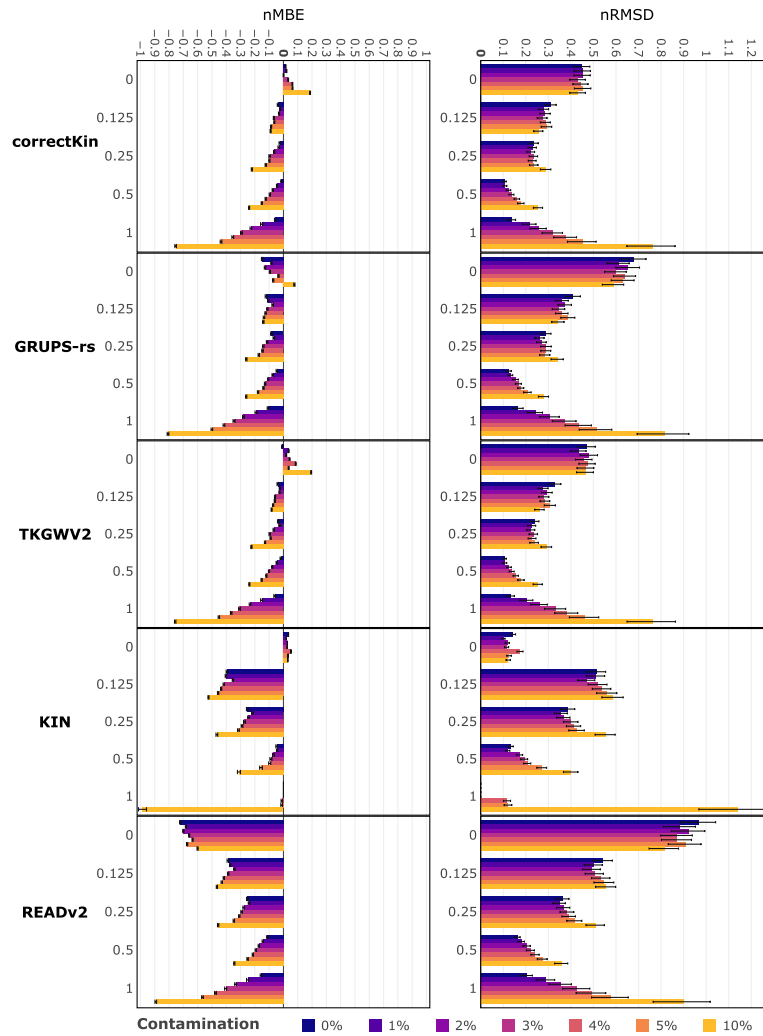


Fig. S8: $nMBE$ (left column) and $nRMSD$ (right column) estimates, across all evaluated methods (rows), rates of contamination from a *GBR* individual (bar colours)[1], and expected relatedness coefficients (y-axis ticks). Increasing values of $nRMSD$ indicate lower accuracy when estimating relatedness coefficients. $nMBE$ values that deviate furthest from zero indicate higher bias, with positive and negative values highlighting a tendency towards over- or under-estimating r-coefficients, respectively. Error bars represent $CI_{95\%}$ for the given estimate.

1.9 Figure S9: Full grid of confusion matrices and UOC values across increasing values of contamination (AFR)

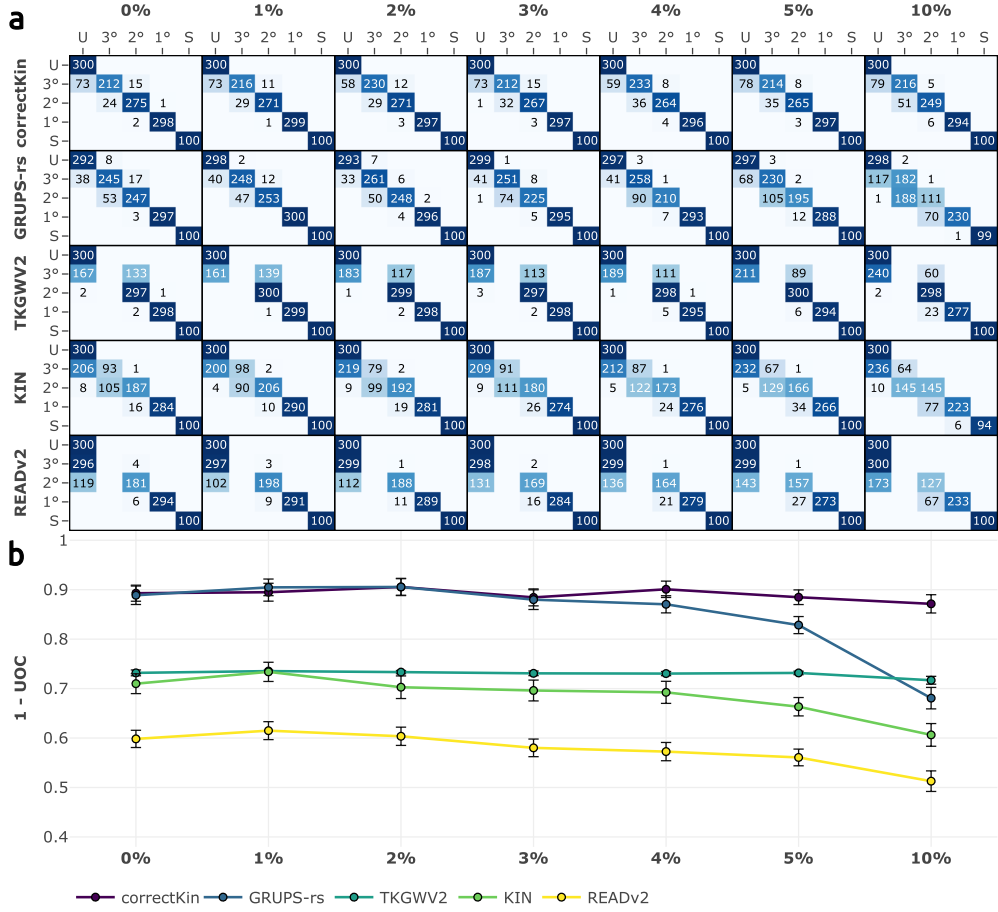


Fig. S9: Benchmark results across increasing values of modern human contamination, using the *AFR* population [1] as a source of contaminating individuals. **a:** Confusion matrices of the five tested methods, confronting expected and predicted relationships. Expected and predicted values are displayed in rows and columns, respectively. 1°, 2°, 3° correspond to first-, second-, and third-degree relationships, respectively, *U* corresponds to "unrelated individuals", and *S* to "self" (monozygotic twins). **b:** UOC values summarizing the classification performance of each method for the considered contamination rate. Higher values of $1 - UOC$ indicate higher performance.

1.10 Figure S10: Full grid of confusion matrices and UOC values across increasing values of contamination (GBR)

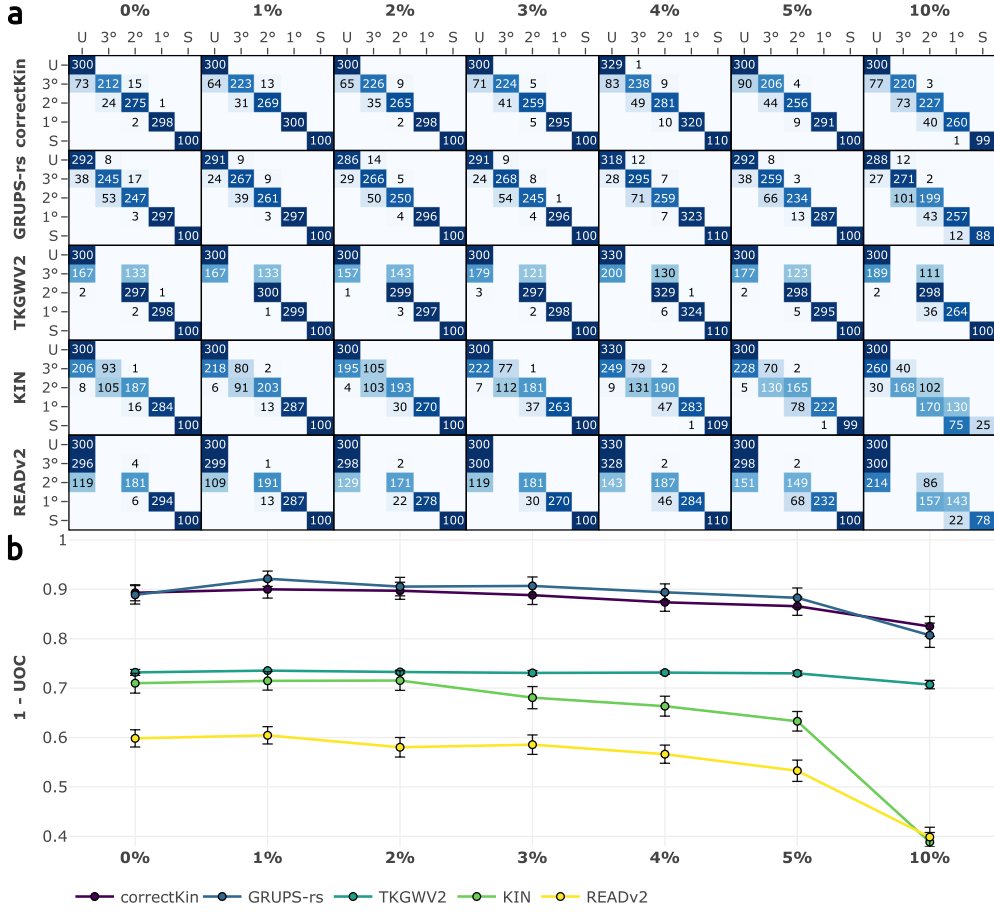


Fig. S10: Benchmark results across increasing values of modern human contamination, using the *GBR* population [1] as a source of contaminating individuals. **a:** Confusion matrices of the five tested methods, confronting expected and predicted relationships. Expected and predicted values are displayed in rows and columns, respectively. 1°, 2°, 3° correspond to first-, second-, and third-degree relationships, respectively, *U* corresponds to "unrelated individuals", and *S* to "self" (monozygotic twins). **b:** UOC values summarizing the classification performance of each method for the considered contamination rate. Higher values of 1 - *UOC* indicate higher performance.

1.11 Figure S11: Impact of admixture on the accuracy and bias of r-coefficients

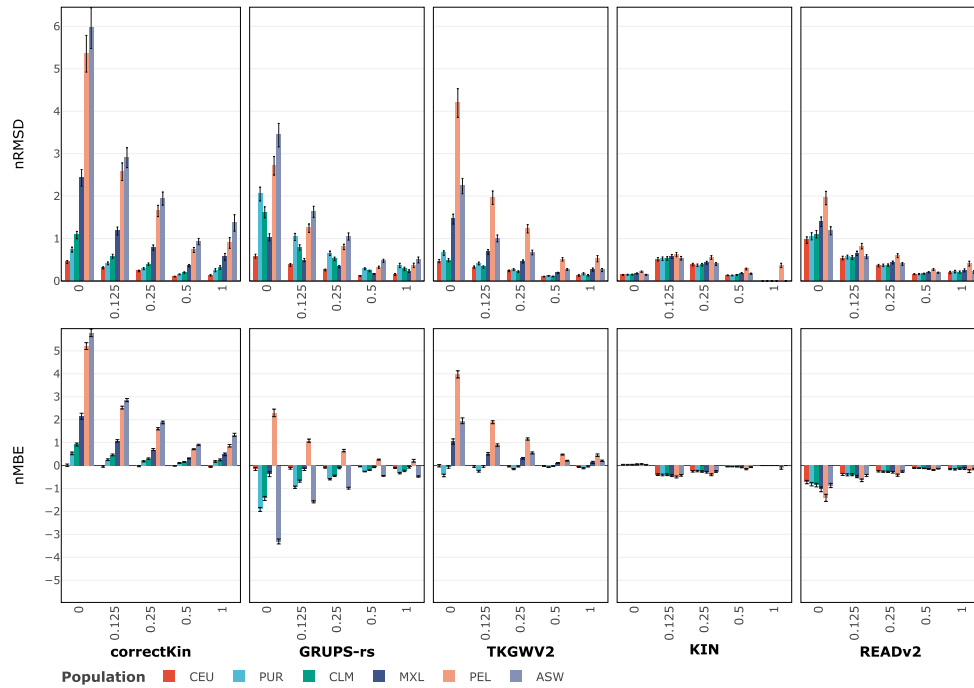


Fig. S11: $nRMSD$ (top row) and $nMBE$ (bottom row) estimates, across all evaluated methods (columns), source populations (bar colours), and expected relatedness coefficients (x-axis ticks), at a simulated sequencing depth of 0.1X. Increasing values of $nRMSD$ indicate lower accuracy when estimating relatedness coefficients. $nMBE$ values that deviate furthest from zero indicate higher bias, with positive and negative values highlighting a tendency towards over- or under-estimating r-coefficients, respectively. Error bars represent $CI_{95\%}$ for the given estimate.

1.12 Figure S12: Confusion matrices and UOC values across increasing values of sequencing depth (ASW)

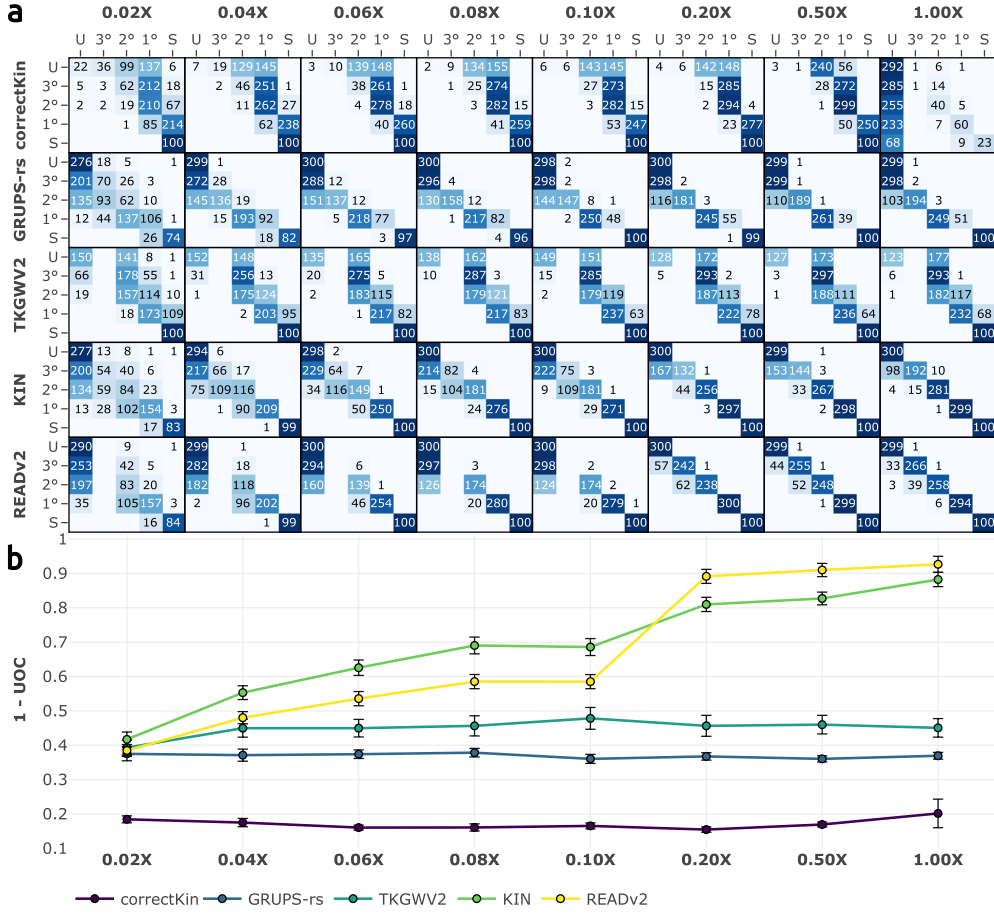


Fig. S12: Benchmark results across increasing values of sequencing depth, using admixed ASW individuals as a source population for pedigree individuals [1]. **a:** Confusion matrices of the five tested methods confronting expected and predicted relationships. Expected and predicted values are displayed in rows and columns, respectively. 1°, 2°, 3° correspond to first-, second-, and third-degree relationships, respectively, U corresponds to "unrelated individuals", and S to "self" (monozygotic twins). **b:** UOC values summarizing the classification performance of each method for the considered sequencing depths. Higher values of $1 - UOC$ indicate higher performance.

1.13 Figure S13: Ancestry proportions of admixed American populations.

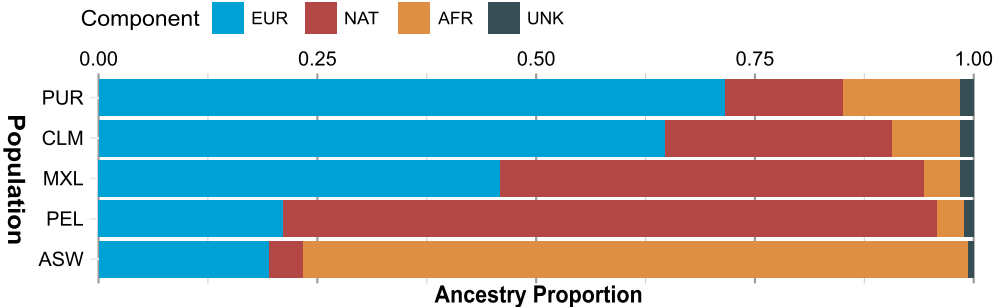


Fig. S13: Ancestry proportions of admixed American populations used during this study. The proportions of this plot were generated using the local ancestry inference results of (Martin et al. 2017) [2] (<https://personal.broadinstitute.org/armartin/tgp-admixture>). **AFR:** African; **EUR:** European; **NAT:** Native American; **UNK:** Unknown

1.14 Figure S14: Average heterozygosity rate of the European
CEU population, and admixed American populations.

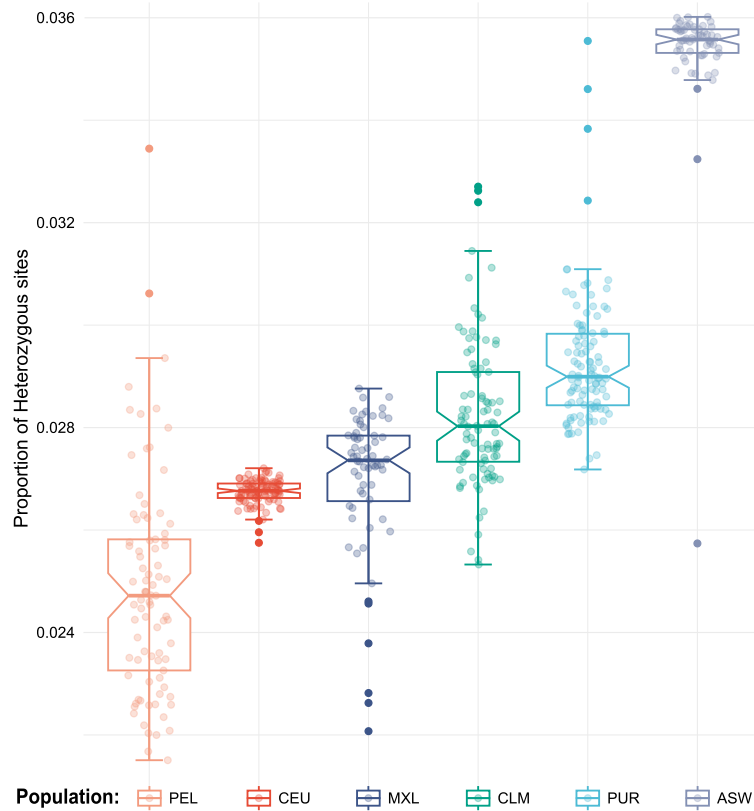
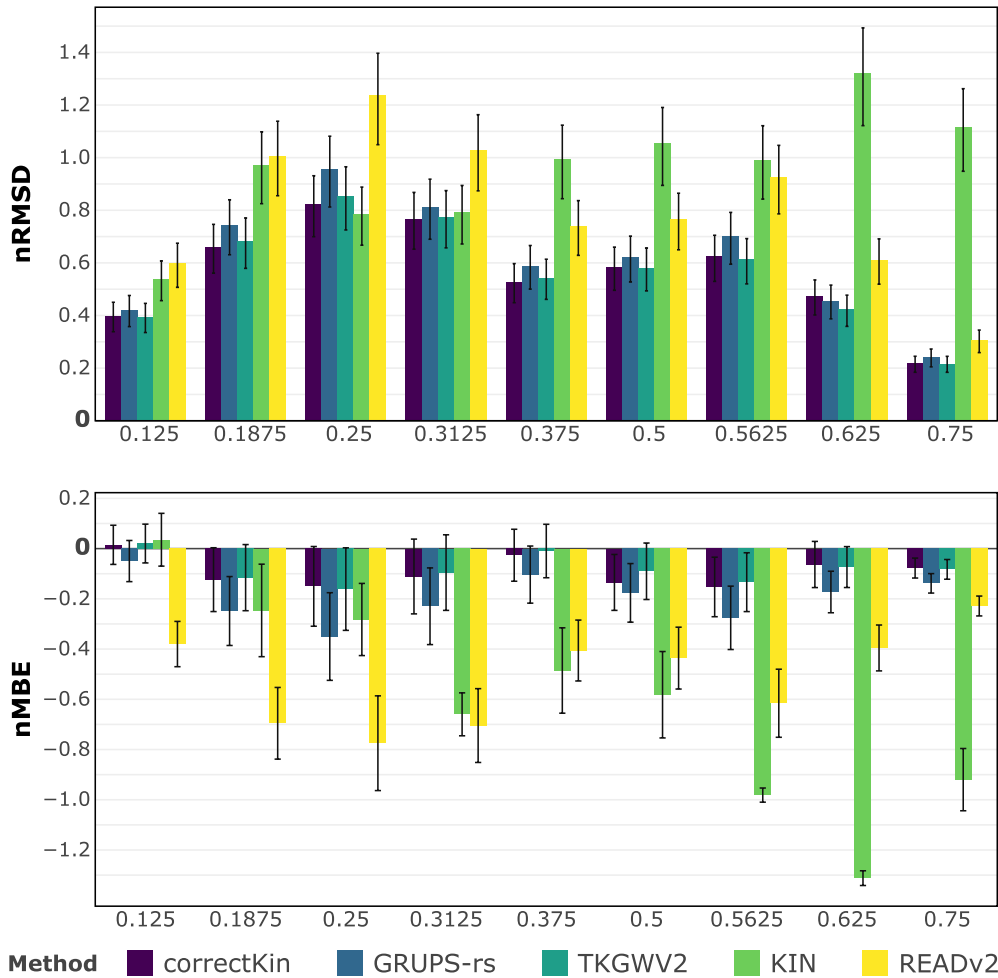


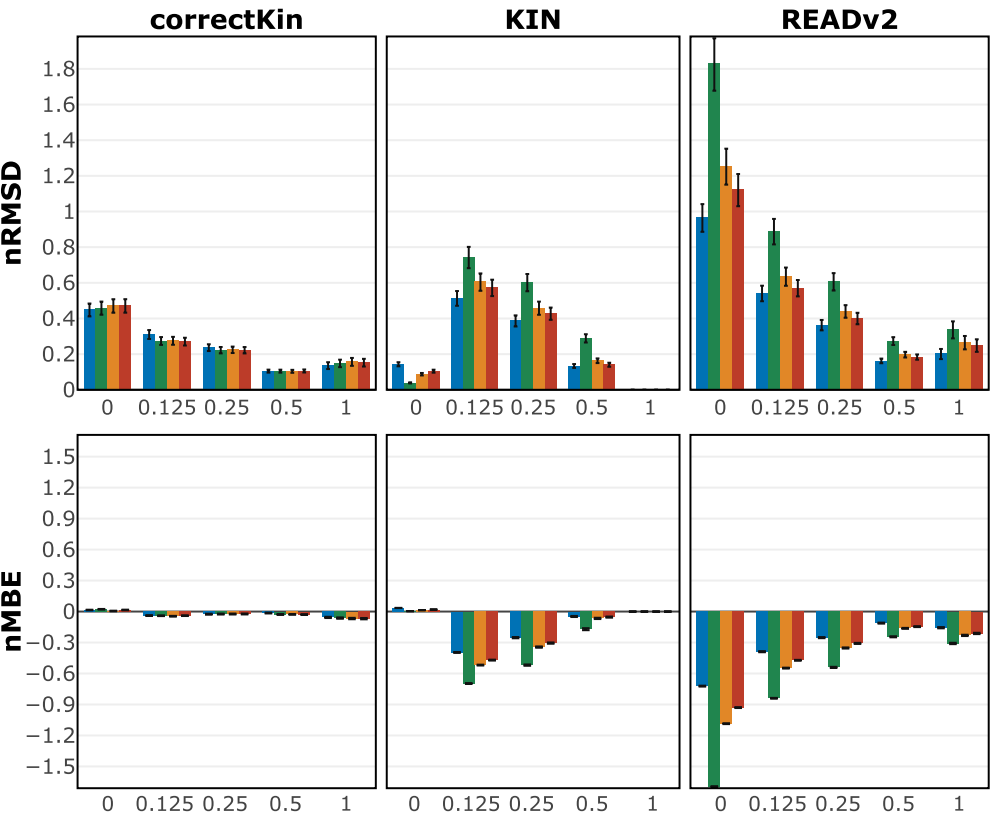
Fig. S14: Sample-wise distribution of the proportions of heterozygous sites of every 1000g-phase3 sample, according to their assigned population. Per-sample counts and proportions of heterozygous sites were directly calculated from the base dataset of the 1000g-phase3 project, using `bcftools stats` [3]. Note that these counts only take SNPs into account. Boxplot notches represent the 95% confidence interval of the median. Whiskers of each boxplot extend from the maximum to the minimum value found within the range $[Q_1 - 1.5 \cdot IQR; Q_3 + 1.5 \cdot IQR]$

1.15 Figure S15: Accuracy and bias of r-coefficients estimates
for pairwise comparisons involving inbred individuals



Method correctKin GRUPS-rs TKGWV2 KIN READv2
Fig. S15: $nRMSD$ (top row) and $nMBE$ (bottom column) estimates, obtained when simulating inbreeding, across all evaluated methods (bar colours) and expected relatedness coefficients (x-axis ticks). Increasing values of $nRMSD$ indicate lower accuracy when estimating relatedness coefficients. $nMBE$ values that deviate furthest from zero indicate higher bias, with positive and negative values highlighting a tendency towards over- or under-estimating r-coefficients, respectively. Error bars represent $CI_{95\%}$ for the given estimate.

1.16 Figure S16: Impact of inbreeding on the accuracy and bias of r-coefficients estimates for pairwise comparisons involving outbred individuals



Pedigree ■ outbred ■ first-cousins ■ half-siblings ■ full-siblings
Fig. S16: *nRMSD* (top row) and *nMBE* (bottom column) estimates, obtained when simulating inbreeding, across all evaluated methods (columns), expected relatedness coefficients (x-axis ticks) and pedigree scenarii (bar colours). Increasing values of *nRMSD* indicate lower accuracy when estimating relatedness coefficients. *nMBE* values that deviate furthest from zero indicate higher bias, with positive and negative values highlighting a tendency towards over- or under-estimating r-coefficients, respectively. Error bars represent $CI_{95\%}$ for the given estimate.

2 Material and Methods

2.1 Description of BADGER's simulation pipeline

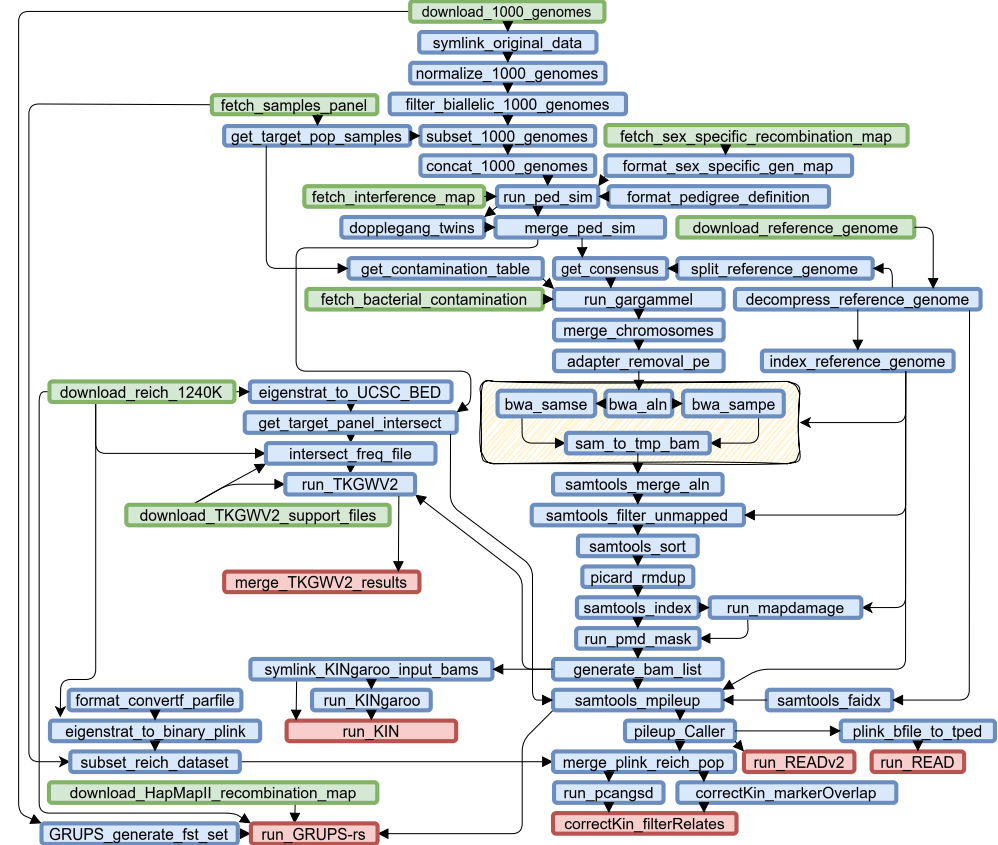


Fig. S17: Complete directed acyclic graph of the BADGER workflow, given the input and parameters provided throughout this study. Node names follow those of the snake-make rules found in BADGER's source code. Green nodes represent data entry points which are automatically downloaded by BADGER. Red nodes represent snakemake rules that are targeted by BADGER by default. The yellow cluster denotes a grouped data input for the reference genome.

2.1.1 1000 genomes dataset pre-processing	921
	922
BADGER first downloads the 1000 genomes phase3-v5b dataset from the EMBL-EBI	923
FTP website [1] (See: Key Resources Table) and proceeds to apply normalization and	924
left-alignment of indels, using <code>bcftools norm</code> . Multi-allelic positions and any lingering	925
unphased genotypes are then filtered out using <code>bcftools view (--phased -m2 -M2)</code> .	926
From the processed dataset, BADGER then generates two data subsets according to	927
the population or super-population label of the samples:	928
	929
• A concatenated VCF file containing all samples belonging to the selected founder	930
population. This file is used as an input to <code>ped-sim</code> when simulating pedigrees, as	931
a source of founder individuals. During this study, the selected founder population	932
was either CEU, to simulate a genetically homogeneous population, or one of the five	933
populations belonging to the admixed American population (ASW, CLM, MXL, PEL,	934
PUR) to simulate an admixed population.	935
	936
• A set of VCF files (one for each autosome), containing all samples belonging to	937
the selected contaminating population is generated. This second dataset is used as	938
input for <code>gargammel</code> when simulating ancient DNA fragments to extract a single	939
contaminating individual, and use its genotype as a source of modern human con-	940
tamination. During this study, we selected either the AFR super-population, or the	941
GBR population as a source to simulate modern human contamination.	942
	943
	944
	945
	946
	947
	948
	949
	950
	951
	952
	953
	954
2.1.2 Pedigree simulations	955
	956
BADGER leverages the software <code>ped-sim</code> to simulate pedigrees in multiple repli-	957
cates, using founder individuals randomly selected from the <code>founder dataset</code> . Here,	958
we parametrized <code>ped-sim</code> to simulate sex-specific recombination rates, as well as a	959
crossover interference model, using the refined genetic map from [4] and the inter-	960
ference parameter estimates of [5], respectively. To maximize the number of possible	961
combinations, and given that BADGER only simulates autosomes, the original genetic	962
	963
	964
	965
	966

967 sex of the individuals selected as founders was not taken into account to select
 968 founders within the pedigree replicates. Simulation of genotyping errors, opposite
 969 homozygous errors, missingness, and pseudo-haploid rates were all disabled at this
 970 step, to prevent any compounding interactions with the error model of `gargammel`
 971 (`--err_rate 0 --err_hom_rate 0 --miss_rate 0 --pseudo_hap 0`).
 972 Simulation of monozygotic twins and/or duplicate samples within the pedigree is
 973 performed by merely duplicating the genotype of the selected individuals within the
 974 output VCF of `ped-sim` (See rule "*dopplegang_twins*", Supplementary Figure S17).
 975
 976
 977
 978
 979
 980
 981 **2.1.3 Ancient DNA simulations**
 982
 983 Simulation of raw ancient DNA fragments for every pedigree individual is performed
 984 using `gargammel`. As this software requires the use of FASTA-format haplotype
 985 sequences, BADGER first uses `bcftools consensus` to apply the variants from the
 986 output VCF file of `ped-sim` to a reference `.fasta` file, thus generating a consen-
 987 sus sequence for every pedigree individual and, when simulating non-null rates of
 988 modern human contamination, a randomly selected contaminating sample from the
 989 `contaminant dataset`.
 990
 991 Here, note that copy-number variations, two-sided inversions, and insertions of
 992 ALU, LINE1, SVA and Nuclear Mitochondrial elements are filtered out using regular
 993 expressions, to comply with the requirements of `bcftools` (`--exclude 'ALT~"<CN`
 994 [`0-9`] `.*>" | ALT~"<INS:.*>" || ALT~"<INV>'`).
 995
 996 For every individual, haplotype sequences are then inserted in the required `endo`
 997 input directory of `gargammel`. Likewise, when simulating human contamination, the
 998 haplotype sequences of the randomly sampled individual are inserted in the optional
 999 `cont` input directory. BADGER then applies `gargammel` on these input directories,
 1000 using the user-provided misincorporation probability and fragment size frequency pro-
 1001 files. Here, we elected to use the *post-mortem* damage profile of "*Chan_Meso*": a young
 1002 adult female individual dated from the Mesolithic period (9137 ± 124 *cal.BP*) and

exhumed from the "Chan do Lindero" karst system of Pedrafita, Lugo, Spain [6]. This choice of reference was motivated by the fact that Chan_Meso was sequenced on an Illumina HiSeq2000 platform – one of the preset sequencing platform model choice for `gargammel`'s – and exhibits "average" *post-mortem* damage patterns (i.e. an approximate misincorporation rate of 0.22, at the 3'-and 5'-end of reads, and a mode of approximately 70 base pairs on its fragment size frequency distribution). Note that while BADGER can be parametrized to handle bacterial contamination from publicly available databases, this capacity was not leveraged during the present study.

To optimize the I/O throughput and runtime performance of BADGER, generation of ancient DNA fragments using `gargammel` is applied in parallel for every pedigree individual, on a per-chromosome basis using a simple scatter-gather approach. Hence, per-chromosome FASTQ files are merely concatenated using `zcat` and `gzip` UNIX command-line utilities.

2.1.4 Alignment

The raw paired-end fragments of every individual composing the pedigree are then trimmed of adapter sequences and collapsed, using `AdapterRemovalv2` [7], requiring a minimum adapter overlap of 1, read length of 17 and base quality of 20. (`--minlength 17 --minquality 20 --minadapteroverlap 1`).

Trimmed fragments are then aligned against the GRCh37 reference genome, using `bwa aln` [8], following the best practices described in [9] (`-l 1024 -n 0.01 -k 2 -o 2`). Note that collapsed single-end sequences and non-collapsed paired-end sequences are mapped separately, using `bwa samse` and `bwa sampe` respectively, and then merged using `samtools merge`. Here, a generic read group tag is placed using `samtools addreplacerg` following merging.

1059 **2.1.5 Quality filtering and preprocessing of alignment files**

1060
1061 Following alignment and merging, a simple quality filtration step is first applied to
1062 the raw binary alignment files of every sample using `samtools view`. Hence, the raw
1064 files are trimmed of any sequence that is either a) unmapped, b) measuring less than
1065 30 nucleotides, or c) carrying a mapping quality score lower than 20 (PHRED scale)
1067 `(-F4 -q20 -e 'length(seq)>30')`.

1069 Alignment files are then sorted using `samtools sort` and removed of any
1070 optical PCR duplicates using `picard MarkDuplicates (--REMOVE_DUPLICATES true`
1071 `--VALIDATION_STRINGENCY LENIENT --ASSUME_SORT_ORDER coordinate)`
1072

1076 **2.1.6 Correction of *post-mortem* deaminations**

1077
1078 Patterns of *post-mortem* deamination were estimated on every sample alignment file,
1079 using `mapDamage2`. To estimate the performance impact of applying *post-mortem*
1081 damage rescaling, two alternative post-processing methods were then applied:
1082

- 1083 • "Rescaled" versions of the alignment files, wherein the base quality scores of putative
1084 nucleotide misincorporation sites are downscaled, were generated by applying the
1086 `--rescale` flag of `mapDamage2`.
1087
- 1088 • "Masked" version of the alignment files were generated using the in-house soft-
1089 ware `pmd-mask` and the misincorporation probability estimates of `mapDamage2`
1091 (`misincorporation.txt` file). Here, potential C>T and G>A deamination sites are
1092 masked all together, along the 5' and 3' ends of fragments, respectively, until the
1093 misincorporation probability is less than 1%.
1094
1095
1096

1098 **2.1.7 Variant calling**

1099
1100 Next, BADGER jointly applies random pseudo-haploid variant calling on every post-
1101 processed alignment file by first creating a pileup file with `samtools mpileup`. Here,
1102 autosomal bi-allelic SNP positions from the AADR "1240K" SNP dataset, version 52.2
1104

were targeted [10], while disabling Base alignment quality (BAQ) recalculation, and filtering out any position with a mapping and/or base quality lower than 20 (`-RB -q20 -Q20`). This pileup file is then directly given as input to the `pileupCaller` module of `sequenceTools` (<https://github.com/stschiff/sequenceTools>), to generate pseudo-haploid variant calls (`--randomHaploid --minDepth 1`), in binary PLINK format.

2.1.8 Genetic relatedness estimation

Note that the benchmarked genetic relatedness estimation methods may have differing input data. Hence:

- The random pseudo-haploid variant calls of `pileupCaller` were given as input to `READv1` and `READv2`.
- The joint pileup file of `samtools mpileup` was given as input to `GRUPS-rs`
- The post-processed binary alignment files of every sample composing a pedigree replicate were given as input to `correctKin`, `KIN` and `TKGWV2`.

correctKin

Following the guidelines of [11], `BADGER` first generates a subset of the AADR "1240K", to provide `correctKin` with a user-selected set of reference individual genotypes. Here, all samples belonging to the EUR super-population of the 1000g-phase3 dataset and contained within the 1240K dataset were selected as reference individuals for `correctKin` during this study. However, as `BADGER` also makes use of 1000-genomes samples as a source of founder individuals during pedigree simulations, the pipeline first excludes any sample previously given as an input to `ped-sim`, from the list of reference samples added to the `correctKin` input dataset. `BADGER` then merges the resulting "1240K" data subset with the pseudo-haploid variant callset of the pedigree replicate, using `plink (--bmerge --merge-mode 1 --allow-no-sex --keep-allele-order)`. Still following guidelines, a covariance matrix and a marker overlap fraction matrix are generated from this merged dataset, using `pcangsd` [12]

1151 and the `markerOverlap` module of `correctKin`. Unrelated individuals were filtered out
1152 using the `filterRelates` module of `correctKin`. Here, note that, a) all pairs of indi-
1153 viduals not found in the resulting output file were considered as unrelated, and b)
1154 pairs of individuals classified as "uncertain" are reclassified as unrelated.
1155

1156

1157

1158 ***GRUPS-rs***

1159

1160 Following the guidelines of [13], `BADGER` first creates an FSA-encoded dataset of
1161 reference individuals from the raw 1000-genomes phase3 dataset, using the `grups-`
1162 `rs fst` module. This preprocessing step is merely intended to increase the runtime
1163 efficiency of `BADGER` and the resulting fsa-encoded 1000g-phase3 dataset is used as
1164 an input throughout all pedigree replicates. For every pedigree replicate, `BADGER`
1165 then directly applies `grups-rs pedigree-sims` on the pileup file described in 2.1.7,
1166 while requesting 1000 simulation replicates, and using samples from the 1000g-phase3
1167 EUR super-population as reference individuals (`--reps 1000 --pedigree-pop EUR`
1168 `--min-depth 1 --seq-error-rate 0.0`).

1169

1170 Since `GRUPS-rs` requires the use of a user-constructed template pedigree to per-
1171 form its simulations, we provided the software with the same template throughout this
1172 study (Supplementary Figure S18). Note that this simple pre-constructed template is
1173 made available as an example within the software's documentation, contains a pair of
1174 siblings, half-siblings and first-cousins, and uses these comparisons to estimate first-,
1175 second- and third-degree relationships, respectively.

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

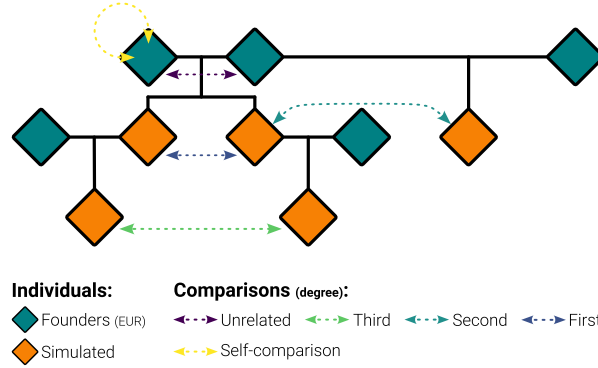


Fig. S18: Diagram of the input template pedigree definition file provided to the GRUPS-rs method throughout this study.

KIN

BADGER first applies the KINGaroo module on the post-processed binary alignment files of all individuals composing a pedigree replicate. Bi-allelic autosomal SNPs from the "1240K" dataset are targeted during this preprocessing step (`--bedfile`), while disabling contamination correction (`--contam_parameter 0`). The main KIN module is then applied on the output of KINGaroo using default parameters. From the final output of KIN, a coefficient of relatedness r is deduced for all tested pairs of individuals, using the provided Cotterman coefficient estimates for every given pair, i.e.: $r = k_1/2 + k_2$.

READv1

The pseudo-haploid variant callset of pileupCaller described in section 2.1.7 is given as input to READ, with default arguments (normalization statistic: median, sliding window size: 10^6 bp).

From the final output of READv1, the relatedness coefficient r of every tested pair of individuals is derived from the normalized \bar{P}_0 estimates of a given pair, using the following equation: $r = 2(1 - \bar{P}_0)$.

1243 **READv2**

1244

1245 Likewise, for every pedigree replicate, the pseudo-haploid variant dataset produced
1246 by `pileupCaller` in section 2.1.7 is provided to READv2 by BADGER, using the
1247 default parameters (`--norm_method median`). Contrary to READv1, the coefficient
1248 of relatedness r is directly obtained from the final output results of the software.
1249

1251

1252

1253 **TKGWV2**

1254

1255 BADGER applies TKGWV2 on every tested pair of post-processed alignment files, using
1256 the support 1000g-phase3 EUR population bed files and allele frequencies provided in
1257 [14] (1000GP3_22M_noFixed_noChr.bed and 1000GP3_EUR_1240K.frq, respectively).
1258 Here, note that providing TKGWV2 with a pre-defined dataset of allele frequencies incurs
1259 the risk of targeting positions that were not simulated by `ped-sim` during a particular
1260 run, as BADGER makes successive use of pedigree simulations using a pre-processed
1261 1000g dataset, followed by the creation of consensus sequences (which will naturally
1262 exclusively contain reference alleles). Hence, to alleviate this potential source of ref-
1263 erence bias, we first filter the provided 1000GP3_EUR_1240K.frq file by removing any
1264 position that was not found within the raw output VCF file previously emitted by `ped`
1265 `-sim`. Also note that, contrary to all other candidate methods, TKGWV2 may only be
1266 applied on a single pair of individuals. Thus, applying this method on an entire pedi-
1267 gree replicate requires BADGER to employ a scatter-gather approach, by running the
1268 software on every tested pair of individuals, and subsequently merging the results.
1269

1278

1279

1280

1281 **2.2 Statistical analysis and benchmark using badger.plots**

1282

1283 Following the application of BADGER in multiple replicates, the statistical analysis
1284 and performance estimation of each method is handled using `badger.plots` : a com-
1285 mand line interface, written as a companion software to BADGER. This software thus
1286 sequentially performs :

1. The deserialization and consolidation of the results of each genetic relatedness estimation software, across all BADGER simulation replicates and sets of studied biological parameters.	1289 1290 1291 1292 1293
2. The calculation of summary statistics regarding the classification performance of each method, for each biological parameter studied.	1294 1295 1296
3. The estimation of the average accuracy and bias of each method's r-coefficient calculation for each degree of relationship.	1297 1298 1299 1300
4. The generation of interactive plots summarizing these performance statistics.	1301 1302 1303 1304
2.2.1 Estimation of classification performance	1305 1306
For every biological condition and method, we constructed confusion matrices confronting the predicted degrees of relationship for all pairwise comparisons, against the "true" degrees of relationship, defined by the topology of the template pedigree originally given as an input for BADGER. From these confusion matrices – one for every pedigree replicate – we calculated the Uniform Ordinal Classification Index (UOC) as a measure of a method's overall classification performance [15]. Briefly, this performance metric, adapted from the ordinal classification index of Cardoso and Sousa [16], is bound between 0 and 1 (0 implying perfectly accurate classification), insensitive to class-imbalance, and markedly takes into account the inherent relative order, and ranking distance separating two degrees of relationship. As such, the ordinal nature of estimating genetic relatedness is retained when estimating performance (e.g. misclassifying an "Unrelated" pairs of individuals as "First-degree" is more penalized than misclassifying them as "Second-degree"). Here, our implementation of the UOC metric was incorporated into <code>badger.plots</code> by adapting the pseudo-code found in [16], and source code provided in [17]. For every method and biological condition, the average UOC of every pedigree replicate is then calculated and plotted as a final aggregate summary statistic. 95% confidence intervals are directly estimated from the	1307 1308 1309 1310 1311 1312 1313 1314 1315 1316 1317 1318 1319 1320 1321 1322 1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334

distribution of UOC values across all simulation replicates, using normal approximation. When applicable, estimates for the area under the curve (AUC) of UOC values of every method were obtained through trapezoidal integration, using the R package `pracma` version 2.4.4, and its associated function `trapz`.

2.2.2 Average accuracy and bias of relatedness coefficients

For every simulated degree of relationship, and across all tested methods and biological conditions, we calculate the average Root Mean Square Deviation (RMSD) and (Mean Bias Error) between the *calculated* and the *expected* relatedness coefficients (r), in an effort to gain insight regarding the average accuracy and bias of each method. Note that, as many of the methods tested here do not *directly* compute r-coefficients, `badger.plots` must first derive this metric from the raw output of every method. Hence:

- As KIN estimates the Jacquard genetic identity coefficients of every pairwise relationship [18], we derived an r-coefficient from the provided k_1 and k_2 values, i.e.:

$$r = \frac{k_1}{2} + k_2$$

- for READ and READv2, an r-coefficient can be calculated from the normalized \overline{P}_0 of a given pair, using the following formula: $r = 2 \cdot (1 - \overline{P}_0)$

- for GRUPS-rs, an r-coefficient can be derived by first calculating normalized estimates of the $PWD_{i,j}^{obs}$ metric of a given pair i, j , which is obtained by dividing this raw estimate by the average expected distribution of unrelated pairs ($\widehat{PWD}_{i,j,unrelated}^{sim}$). It follows that the r-coefficient can be derived using the following equation: $r = 2(1 - \frac{PWD_{i,j}^{obs}}{\widehat{PWD}_{i,j,unrelated}^{sim}})$

- Finally, as both `correctKin` and `TKGWV2` compute a kinship coefficient (ϕ), the r-coefficient is simply obtained by multiplying this estimate by 2: $r = 2 \cdot \phi$

Here, it must be noted that the distance separating the expected r-coefficient of a given degree of relationship from neighbouring distributions varies with the degree of

relatedness, and is effectively halved for each additional degree separating two individuals. This implies that a given deviation from the expected average can have a greatly differing impact on the general accuracy, depending on the degree for which it is observed (e.g. a standard deviation of 0.1 for the r-coefficient between two individuals is insignificant when considering first-degree relationships, but would consistently cause misclassifications in the case of third-degree relationships).

Thus, to properly compare the accuracy and bias, both across a given method and degree of relatedness, we propose to first normalize the RMSD and MBE of a given relationship by the range of its theoretical distribution. This can be done by dividing the RMSD of MBE value by the distance separating the two midpoints found between the mean of a given relationship k and its neighbouring ones ($k - 1$ and $k + 1$).

$$nRMSD_{m,k} = \frac{\sqrt{\frac{\sum_{i=1}^{i=N} (\hat{r}_k - r_{m,k,i})^2}{N}}}{\frac{\hat{r}_{k+1} - \hat{r}_{k-1}}{2}} \quad (1)$$

$$nMBE_{m,k} = \frac{\frac{\sum_{i=1}^{i=N} (r_{m,k,i} - \hat{r}_k)}{N}}{\frac{\hat{r}_{k+1} - \hat{r}_{k-1}}{2}} \quad (2)$$

where m and k represent a given method and degree of relatedness, respectively. \hat{r}_k is the expected relatedness coefficient of relationship k ; $r_{m,k,i}$, the calculated relatedness coefficient for the i^{th} pair of individuals, and N , the total amount of observations of a given relationship k ($k = 0, 0.125, 0.25, 0.5, 1$). 95% confidence intervals for these metrics were calculated using the principles and methods described in [19, 20], i.e.:

$$nRMSD_{m,k} \in \left[nRMSD_{m,k} \left(1 - \sqrt{1 - \frac{1.96\sqrt{2}}{\sqrt{N-1}}} \right); RMSD_{m,k} \left(\sqrt{1 + \frac{1.96\sqrt{2}}{\sqrt{N-1}}} - 1 \right) \right] \quad (3)$$

$$nMBE_{m,k} \in \left[nMBE_{m,k} \pm \frac{1.96 \cdot \sigma_{m,k}}{\sqrt{N}} \right] \quad (4)$$

where $\sigma_{m,k}$ is the population standard deviation of the relatedness coefficients obtained from method m , and belonging to relationship k .

2.3 Key Resources Table

Reagent or Resource	Source	Identifier
Deposited data		
1000g-phase3-v20130502	IGSR [1]	https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/
HapMapII	IHMP [21]	http://ftp.ncbi.nlm.nih.gov/hapmap/recombination/2011-01-phaseII.B37/
GRCh37-release113	Church et al. [22]	http://ftp.ensembl.org/pub/grch37/release-113/fasta/homo_sapiens/dna/
Cross-over interference model	Campbell et al. [5]	https://raw.githubusercontent.com/williamslab/ped-sim/refs/heads/master/interfere/nu_p_campbell.tsv
Sex-specific genetic maps	Bhérier et al. [4]	https://github.com/cbherer/Bherer_etal_SexualDimorphismRecombination/
AADR dataset v52.2	Mallick et al. [10]	https://reichdata.hms.harvard.edu/pub/datasets/amh_repo/curated_releases/V52/V52.2/SHARE/public.dir/
TKGWV2 support files	Fernandes et al. [14]	https://github.com/danimfernandes/tkgwv2
Softwares and algorithms		
AdapterRemoval-v2.3.3	Schubert et al. [7]	RRID:SCR_011834
ANGSD-v0.939	Korneliussen et al. [23]	RRID:SCR_021865
BADGER-v0.5.1	This study	https://github.com/MaelLefevre/badger/tree/v0.5.1
bcftools-1.15	Li [3]	RRID:SCR_005227
conda-23.1.0	Conda contributors [24]	RRID:SCR_018317
correctKin	Nyerki et al. [11]	RRID:SCR_026952
gargammel-1.1.4	Renaud et al. [25]	RRID:SCR_026953

grups-rs-0.3.2	Lefeuvre et al. [13]	RRID:SCR_026954
kin-3.1.3	Popli et al. [26]	RRID:SCR_026955
mapDamage-v2.2.1	Jónsson et al. [27]	RRID:SCR_001240
pcangsd-0.99	Meisner and Albrechtsen [12]	RRID:SCR_026956
ped-sim-v1.4	Caballero et al. [28]	RRID:SCR_026957
picard-v2.27.4	Broad Institute [29]	RRID:SCR_006525
plink-v1.9	Chang et al. [30]	RRID:SCR_001757
pmd-mask-v0.3.2	This study	https://github.com/MaelLefeuvre/pmd-mask/tree/v0.3.2
READ-v1.0	Kuhn et al. [31]	RRID:SCR_026958
READv2-v2.00	Alaçamlı et al. [32]	RRID:SCR_026959
samtools-v1.15	Li [3]	RRID:SCR_002105
sequenceTools-v1.5.2	Schiffels [33]	https://github.com/stschiff/sequenceTools/tree/v1.5.2
snakemake-7.20.0	Mölder et al. [34]	RRID:SCR_003475
TKGWV2	Fernandes et al. [14]	RRID:SCR_026960
python-3.11.0	Python Software Foundation	RRID:SCR_008394
R-v4.1.2	R Development Core Team	RRID:SCR_001905

3 Description of the pmd-mask command line utility

3.1 Rationale, behaviour and workflow description

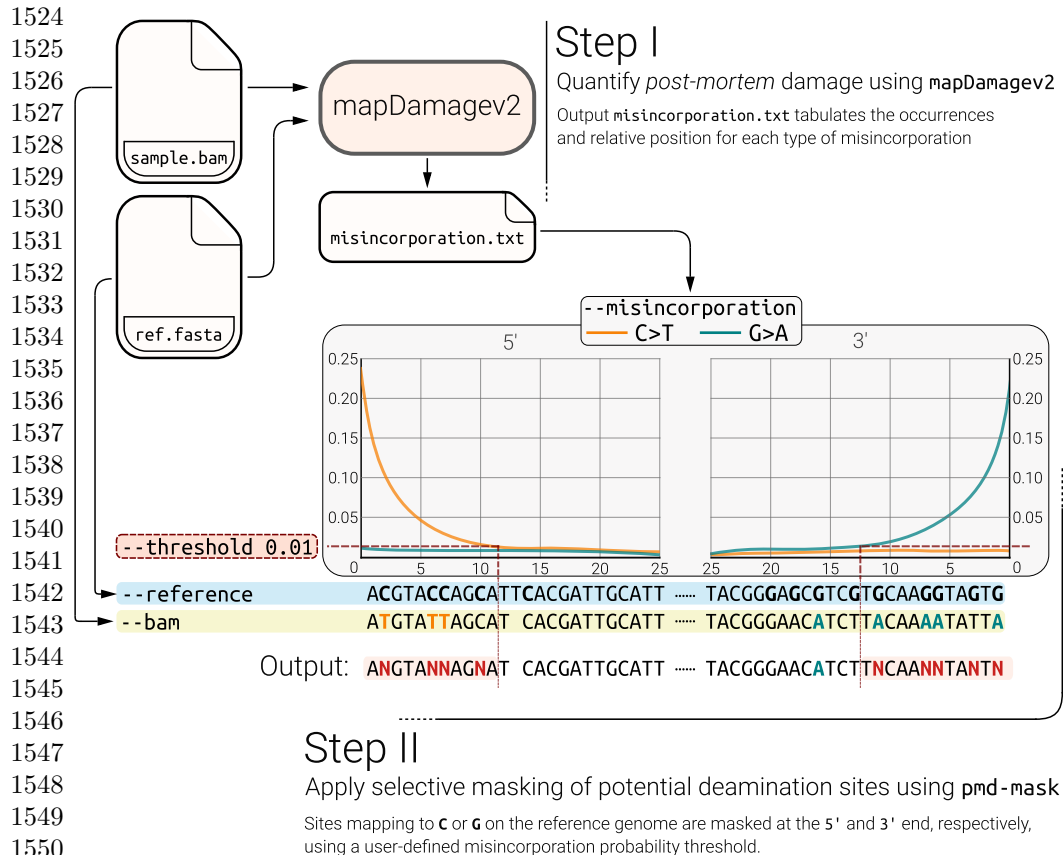


Fig. S19: Summary diagram describing the main process, required inputs, and workflow surrounding the `pmd-mask` command-line utility.

`pmd-mask` is a simple command-line pre-processing tool written in the Rust programming language, which masks the positions from reads that are likely to be impacted by *post-mortem* damage (PMD). Briefly, `pmd-mask` leverages nucleotide- and position-specific misincorporation rate estimates emitted from the `mapDamage2` software (typically, in the form of a `misincorporation.txt` output file) to apply soft-trimming on the read extremities, until the PMD rate reaches a user defined threshold

by setting candidate nucleotides to ‘N’ and their base-quality to 0. Note that the default threshold is here defined as a misincorporation rate of 0.01, but may be modified at leisure by the user, using the `--threshold` argument. Note that, in its current state, the `pmd-mask` algorithm only considers the misincorporation rate found at a given position, and the genotype found in the *reference* genome. Therefore, nucleotides at the extremities of a read are masked, regardless of the actual genotype observed at a given position (Supplementary Figure S19). Pseudo-code snippets, summarizing the main loop of the program can be read in Algorithm-1. Here, the devised approach is one that i) is expected to incur less potential bias than when applying PMD-rescaling through mapDamage2 (i.e. using its provided `--rescale` flag), and ii) carries the benefit of mitigating the loss of information usually displayed when applying hard-clipping, by instead specifically targeting potential C>T and G>A transition sites on both the 5’ and 3’ end of the read, respectively. In other terms, this method may be regarded as a conservative compromise between *post-mortem* damage rescaling methods such as mapDamage2, PMDtools, or ATLAS [27, 35, 36] and hard-clipping methods such as the trimBam module of the bamUtil software [37]. Required inputs for `pmd-mask` are as follows:

- `--bam`: An input `.bam` file (SAM, BAM, and CRAM formats are accepted). `pmd-mask` can either read from a file (using `-b|--bam`) or from the standard input, through shell piping.
- `--misincorporation`: A mapDamage2 `misincorporation.txt` file. This file provides strand-specific PMD frequency estimates, which are used to compute the threshold at which masking should be performed. Evidently, this file must have been generated from the input BAM file to provide a sound estimate.
- `--reference`: A reference genome, in the form of a `.fasta` file. This genome must of course be identical to the one used to align the input BAM file.

1611 Additional instructions regarding the installation and usage of `pmd-mask`, as well
1612 as its source code is made publicly available at [https://github.com/MaelLefevre/](https://github.com/MaelLefevre/pmd-mask)
1613 `pmd-mask`, under GPL-v3.0 licencing.
1614

1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656

3.2 Pseudo-code describing the main algorithm of pmd-mask

Algorithm 1 Pseudo-code describing the main algorithm of pmd-mask.

Ensure: : $args.threshold \in [0.00, 1.00[$

```

1:  $bam \leftarrow bam\_reader(args.bam\_path);$ 
2:  $refseq \leftarrow fasta\_reader(args.reference\_path)$ 
3:  $misincorporation \leftarrow misincorporation\_reader(args.misincorporation\_path)$ 
4:
5: if  $args.threshold$  is not null then
6:    $threshold \leftarrow args.threshold$ 
7: else
8:    $threshold \leftarrow 0.01$ 
9:
10: function MASK(read, position)
11:    $read.nucleotide[position] \leftarrow 'N';$ 
12:    $read.quality[position] \leftarrow 0$ 
13:
14:  $output\_bam \leftarrow bam.copy\_header()$ 
15: for read  $\in bam$  do
16:   ▷ Extract read length, coordinate and strand information
17:    $n \leftarrow read.length$ 
18:    $chr \leftarrow read.chromosome$ 
19:    $pos \leftarrow read.position$ 
20:    $strand \leftarrow read.strand$ 
21:   ▷ Mask 5' Cytosines
22:   for mask5' : ( $i = 0; i < n; i++$ ) do
23:     if  $misincorporation['C > T'][chr][strand][i] \leq args.threshold$  then
24:       break mask5'
25:     else if  $reference.get(chr, pos + i) == 'C'$  then
26:       MASK(read, i)
27:   ▷ Mask 3' Guanines
28:   for mask3' : ( $i = n; i > 0; i--$ ) do
29:     if  $misincorporation['G > A'][chr][strand][i] \leq args.threshold$  then
30:       break mask3'
31:     else if  $reference.get(chr, pos + i) == 'G'$  then
32:       MASK(read, i)
33:
34:   ▷ Store masked read
35:    $output\_bam += read$ 
   return ( $output\_bam$ )

```

References

- [1] The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature*. 2015 Oct;526(7571):68–74. Number: 7571. <https://doi.org/10.1038/nature15393>.
- [2] Martin AR, Gignoux CR, Walters RK, Wojcik GL, Neale BM, Gravel S, et al. Human Demographic History Impacts Genetic Risk Prediction across Diverse Populations. *The American Journal of Human Genetics*. 2017 Apr;100(4):635–649. <https://doi.org/10.1016/j.ajhg.2017.03.004>.
- [3] Li H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*. 2011 Nov;27(21):2987–2993. <https://doi.org/10.1093/bioinformatics/btr509>.
- [4] Bhérier C, Campbell CL, Auton A. Refined genetic maps reveal sexual dimorphism in human meiotic recombination at multiple scales. *Nature Communications*. 2017 Apr;8(1):14994. Number: 1. <https://doi.org/10.1038/ncomms14994>.
- [5] Campbell CL, Furlotte NA, Eriksson N, Hinds D, Auton A. Escape from crossover interference increases with maternal age. *Nature Communications*. 2015 Feb;6(1):6260. Number: 1. <https://doi.org/10.1038/ncomms7260>.
- [6] González-Fortes G, Jones ER, Lightfoot E, Bonsall C, Lazar C, Grandal-d’Anglade A, et al. Paleogenomic Evidence for Multi-generational Mixing between Neolithic Farmers and Mesolithic Hunter-Gatherers in the Lower Danube Basin. *Current Biology*. 2017 Jun;27(12):1801–1810.e10. <https://doi.org/10.1016/j.cub.2017.05.023>.

- [7] Schubert M, Lindgreen S, Orlando L. AdapterRemoval v2: rapid adapter trimming, identification, and read merging. BMC Research Notes. 2016 Feb;9(1):88. <https://doi.org/10.1186/s13104-016-1900-2>.
- [8] Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics (Oxford, England). 2009 Jul;25(14):1754–1760. Number: 14. <https://doi.org/10.1093/bioinformatics/btp324>.
- [9] Oliva A, Tobler R, Cooper A, Llamas B, Souilmi Y. Systematic benchmark of ancient DNA read mapping. Briefings in Bioinformatics. 2021 Sep;22(5):bbab076. <https://doi.org/10.1093/bib/bbab076>.
- [10] Mallick S, Micco A, Mah M, Ringbauer H, Lazaridis I, Olalde I, et al. The Allen Ancient DNA Resource (AADR) a curated compendium of ancient human genomes. Scientific Data. 2024 Feb;11(1):182. <https://doi.org/10.1038/s41597-024-03031-7>.
- [11] Nyerki E, Kalmár T, Schütz O, Lima RM, Neparáczki E, Török T, et al. correctKin: an optimized method to infer relatedness up to the 4th degree from low-coverage ancient human genomes. Genome Biology. 2023 Feb;24(1):38. Number: 1. <https://doi.org/10.1186/s13059-023-02882-4>.
- [12] Meisner J, Albrechtsen A. Inferring Population Structure and Admixture Proportions in Low-Depth NGS Data. Genetics. 2018 Oct;210(2):719–731. <https://doi.org/10.1534/genetics.118.301336>.
- [13] Lefeuvre M, Martin MD, Jay F, Marsolier MC, Bon C. GRUPS-rs, a high-performance ancient DNA genetic relatedness estimation software relying on pedigree simulations. Human Population Genetics and Genomics. 2024 Jan;4(1). <https://doi.org/10.47248/hpgg2404010001>.

1795 [14] Fernandes DM, Cheronet O, Gelabert P, Pinhasi R. TKGWV2: an ancient
1796 DNA relatedness pipeline for ultra-low coverage whole genome shotgun data.
1797 Scientific Reports. 2021 Oct;11(1):21262. Number: 1. [https://doi.org/10.1038/](https://doi.org/10.1038/s41598-021-00581-3)
1798 [s41598-021-00581-3](https://doi.org/10.1038/s41598-021-00581-3).
1800
1801
1802 [15] Silva W, Pinto JR, Cardoso JS. A Uniform Performance Index for Ordinal Clas-
1803 sification with Imbalanced Classes. In: 2018 International Joint Conference on
1804 Neural Networks (IJCNN); 2018. p. 1–8. ISSN: 2161-4407.
1805
1806
1807
1808 [16] Cardoso JS, Sousa R. Measuring the performance of ordinal classification.
1809 International Journal of Pattern Recognition and Artificial Intelligence. 2011
1810 Dec;25(08):1173–1195. <https://doi.org/10.1142/S0218001411009093>.
1811
1812
1813
1814 [17] Albuquerque T, Cruz R, Cardoso JS. Ordinal losses for classification of cervical
1815 cancer risk. PeerJ Computer Science. 2021 Apr;7:e457. [https://doi.org/10.7717/](https://doi.org/10.7717/peerj-cs.457)
1816 [peerj-cs.457](https://doi.org/10.7717/peerj-cs.457).
1817
1818
1819
1820 [18] Jacquard A. Genetic Information Given by a Relative. Biometrics.
1821 1972;28(4):1101–1114. <https://doi.org/10.2307/2528643>.
1822
1823
1824 [19] Nicholls A. Confidence limits, error bars and method comparison in molec-
1825 ular modeling. Part 1: The calculation of confidence intervals. Journal of
1826 Computer-Aided Molecular Design. 2014 Sep;28(9):887–918. [https://doi.org/10.](https://doi.org/10.1007/s10822-014-9753-z)
1827 [1007/s10822-014-9753-z](https://doi.org/10.1007/s10822-014-9753-z).
1828
1829
1830
1831 [20] Nicholls A. Confidence limits, error bars and method comparison in molecular
1832 modeling. Part 2: comparing methods. Journal of Computer-Aided Molecular
1833 Design. 2016 Feb;30(2):103–126. <https://doi.org/10.1007/s10822-016-9904-5>.
1834
1835
1836
1837 [21] Frazer KA, Ballinger DG, Cox DR, Hinds DA, Stuve LL, Gibbs RA, et al. A
1838 second generation human haplotype map of over 3.1 million SNPs. Nature. 2007
1839
1840

- Oct;449(7164):851–861. <https://doi.org/10.1038/nature06258>. 1841
- [22] Church DM, Schneider VA, Graves T, Auger K, Cunningham F, Bouk N, 1842
et al. Modernizing Reference Genome Assemblies. PLOS Biology. 2011 1843
Jul;9(7):e1001091. <https://doi.org/10.1371/journal.pbio.1001091>. 1844
1845
1846
1847
1848
- [23] Korneliussen TS, Albrechtsen A, Nielsen R. ANGSD: Analysis of Next Generation 1849
Sequencing Data. BMC Bioinformatics. 2014 Nov;15(1):356. <https://doi.org/10.1186/s12859-014-0356-4>. 1850
1851
1852
1853
1854
- [24] Conda contributors.: conda: A system-level, binary package and environment 1855
manager running on all major operating systems and platforms. Available from: 1856
<https://docs.conda.io/projects/conda/>. 1857
1858
1859
1860
- [25] Renaud G, Hanghøj K, Willerslev E, Orlando L. gargammel: a sequence simulator 1861
for ancient DNA. Bioinformatics. 2017 Feb;33(4):577–579. Number: 4. <https://doi.org/10.1093/bioinformatics/btw670>. 1862
1863
1864
1865
1866
- [26] Popli D, Peyrégne S, Peter BM. KIN: a method to infer relatedness from low- 1867
coverage ancient DNA. Genome Biology. 2023 Jan;24(1):10. Number: 1. <https://doi.org/10.1186/s13059-023-02847-7>. 1868
1869
1870
1871
1872
- [27] Jónsson H, Ginolhac A, Schubert M, Johnson PLF, Orlando L. mapDam- 1873
age2.0: fast approximate Bayesian estimates of ancient DNA damage parameters. 1874
Bioinformatics. 2013 Jul;29(13):1682–1684. Number: 13. <https://doi.org/10.1093/bioinformatics/btt193>. 1875
1876
1877
1878
1879
- [28] Caballero M, Seidman DN, Qiao Y, Sannerud J, Dyer TD, Lehman DM, et al. 1880
Crossover interference and sex-specific genetic maps shape identical by descent 1881
sharing in close relatives. PLOS Genetics. 2019 Dec;15(12):e1007979. Number: 1882
12. <https://doi.org/10.1371/journal.pgen.1007979>. 1883
1884
1885
1886

1887 [29] Broad Institute.: Picard: A set of command line tools (in Java) for manipulating
1888 high-throughput sequencing (HTS) data and formats such as SAM/BAM/CRAM
1889 and VCF. Broad Institute. Available from: [http://broadinstitute.github.io/](http://broadinstitute.github.io/picard)
1890 [picard](http://broadinstitute.github.io/picard).
1891
1892
1893
1894 [30] Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ.
1895 Second-generation PLINK: rising to the challenge of larger and richer
1896 datasets. GigaScience. 2015 Dec;4(1):s13742–015–0047–8. [https://doi.org/10.](https://doi.org/10.1186/s13742-015-0047-8)
1897 [1186/s13742-015-0047-8](https://doi.org/10.1186/s13742-015-0047-8).
1898
1899
1900
1901
1902 [31] Kuhn JMM, Jakobsson M, Günther T. Estimating genetic kin relationships in
1903 prehistoric populations. PLOS ONE. 2018 Apr;13(4):e0195491. Number: 4. [https:](https://doi.org/10.1371/journal.pone.0195491)
1904 [//doi.org/10.1371/journal.pone.0195491](https://doi.org/10.1371/journal.pone.0195491).
1905
1906
1907
1908 [32] Alaçamlı E, Naidoo T, Güler MN, Sağlıcan E, Aktürk S, Mapelli I, et al.
1909 READv2: advanced and user-friendly detection of biological relatedness in
1910 archaeogenomics. Genome Biology. 2024 Aug;25(1):216. [https://doi.org/10.1186/](https://doi.org/10.1186/s13059-024-03350-3)
1911 [s13059-024-03350-3](https://doi.org/10.1186/s13059-024-03350-3).
1912
1913
1914
1915 [33] Schiffels S.: sequenceTools. Available from: [https://github.com/stschiff/](https://github.com/stschiff/sequenceTools)
1916 [sequenceTools](https://github.com/stschiff/sequenceTools).
1917
1918
1919 [34] Mölder F, Jablonski KP, Letcher B, Hall MB, Tomkins-Tinch CH, Sochat V,
1920 et al.: Sustainable data analysis with Snakemake. F1000Research. Available from:
1921 <https://f1000research.com/articles/10-33>.
1922
1923
1924
1925 [35] Skoglund P, Northoff BH, Shunkov MV, Derevianko AP, Pääbo S, Krause J,
1926 et al. Separating endogenous ancient DNA from modern day contamination in
1927 a Siberian Neandertal. Proceedings of the National Academy of Sciences. 2014
1928 Feb;111(6):2229–2234. <https://doi.org/10.1073/pnas.1318934111>.
1929
1930
1931
1932

[36] Link V, Kousathanas A, Veeramah K, Sell C, Scheu A, Wegmann D.: ATLAS: Analysis Tools for Low-depth and Ancient Samples. bioRxiv. Pages: 105346 Section: New Results. Available from: <https://www.biorxiv.org/content/10.1101/105346v2>.

[37] Jun G, Wing MK, Abecasis GR, Kang HM. An efficient and scalable analysis framework for variant extraction and refinement from population scale DNA sequence data. Genome Research. 2015 Apr;p. gr.176552.114. Company: Cold Spring Harbor Laboratory Press Distributor: Cold Spring Harbor Laboratory Press Institution: Cold Spring Harbor Laboratory Press Label: Cold Spring Harbor Laboratory Press Publisher: Cold Spring Harbor Lab. <https://doi.org/10.1101/gr.176552.114>.