

aRchaic

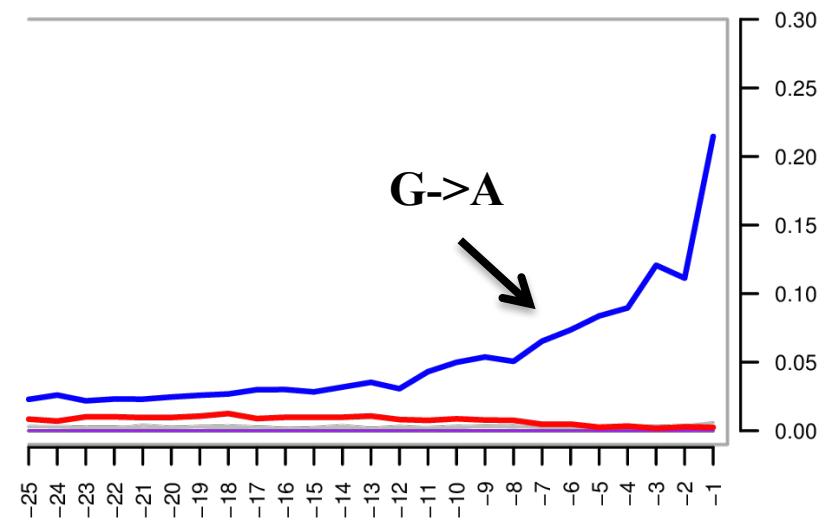
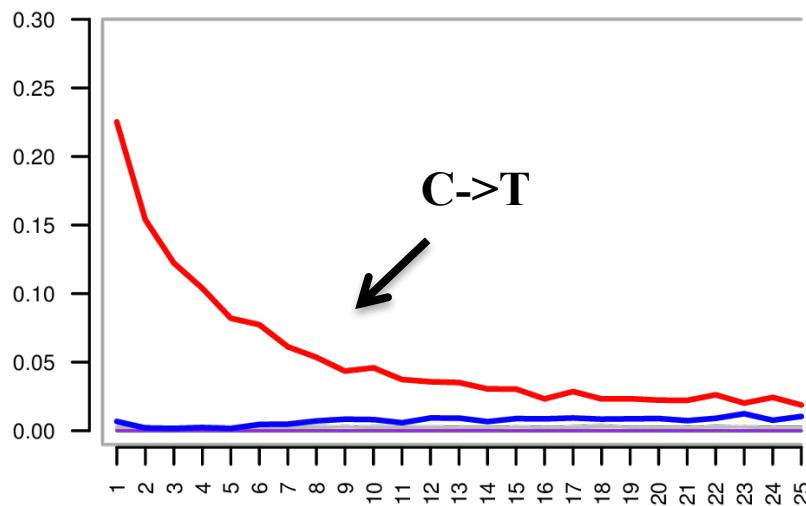


An R package for clustering and
visualizing ancient dna signatures

By Hussein Al-Asadi & Kushal Dey

mapDamage

mapDamage: tracking and quantifying
damage patterns in ancient DNA sequences

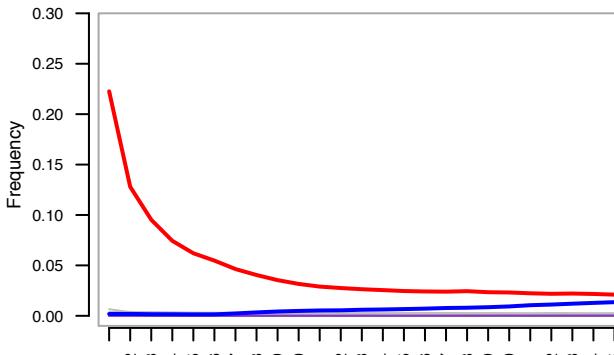


Data-set # 1

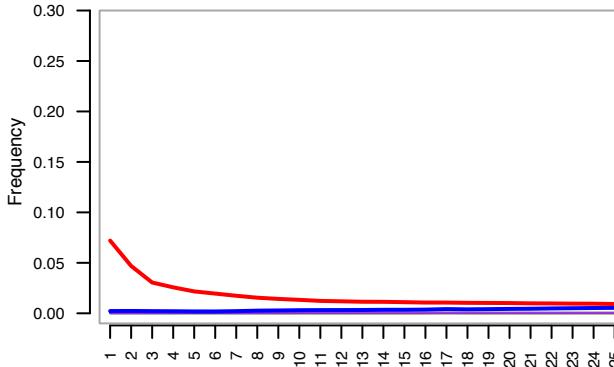
- DiRenzo Lab: 48 Ancients and 7 control samples
- Questions:
 - Do the ancients look “ancient”?
 - Is there any DNA in the controls?
 - If so, is the DNA modern or ancient?

MapDamage on Ancients

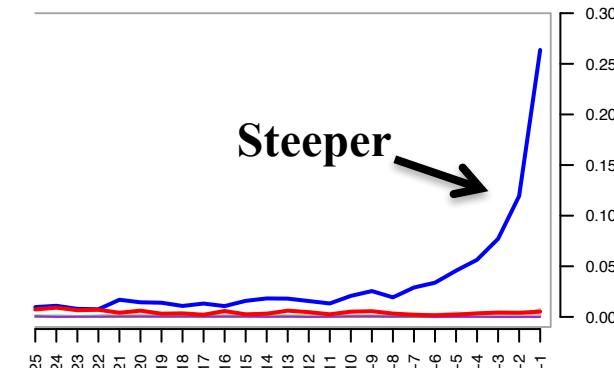
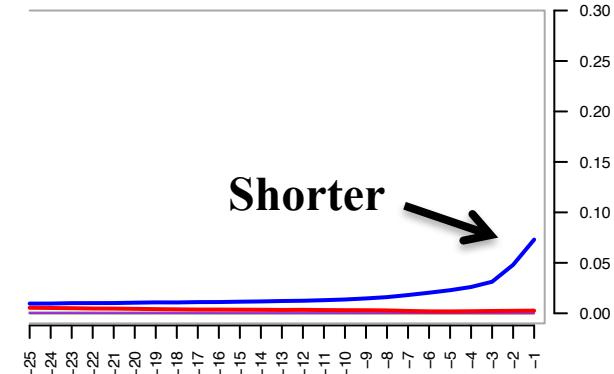
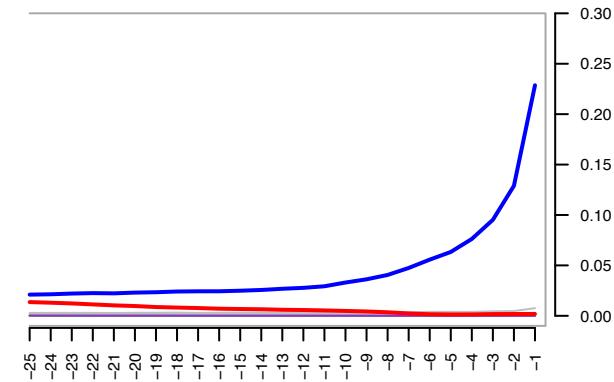
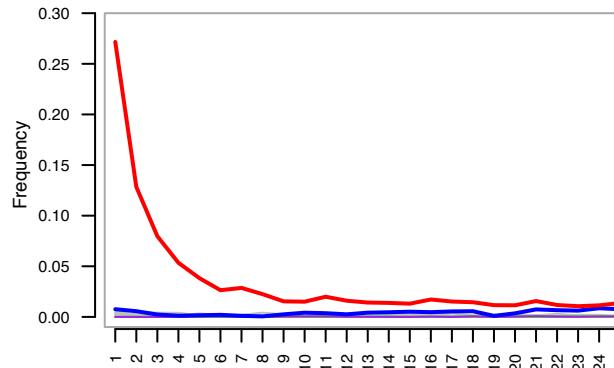
CNE1



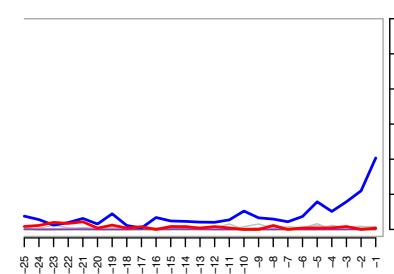
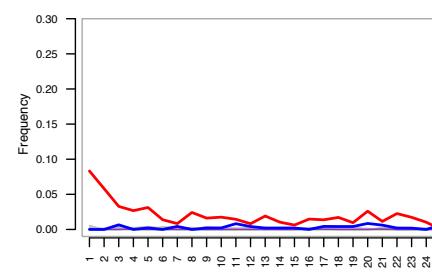
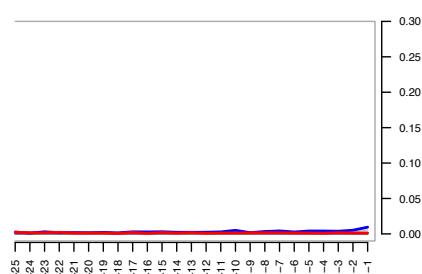
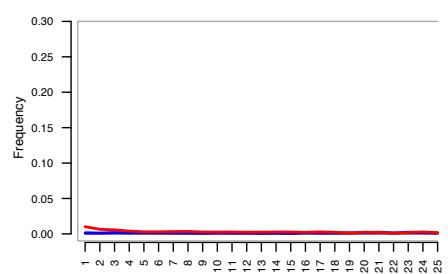
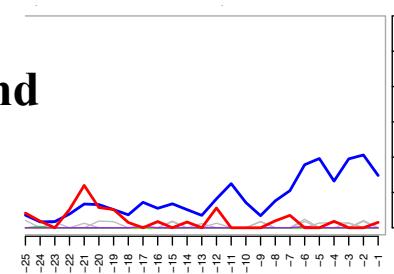
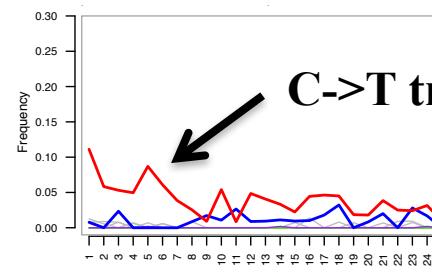
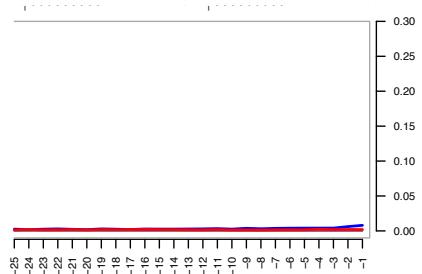
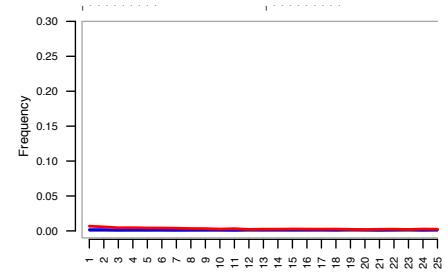
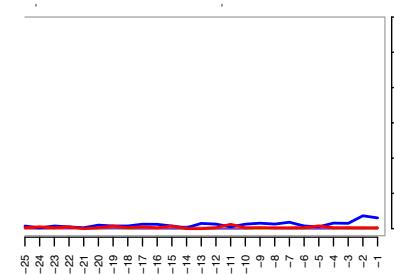
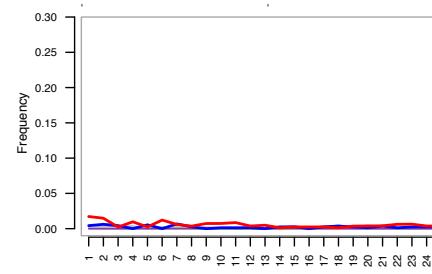
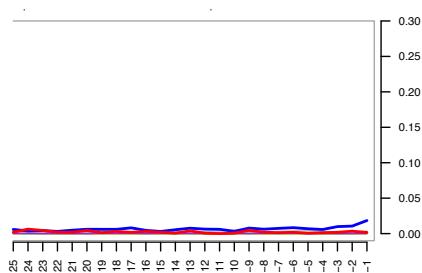
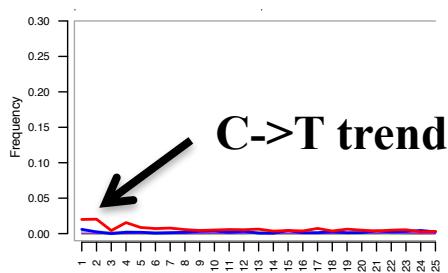
KS20



S30

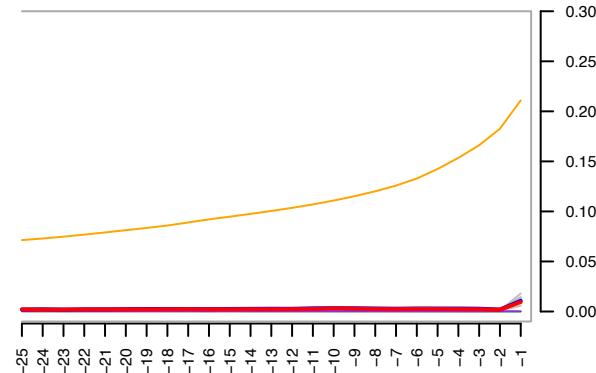
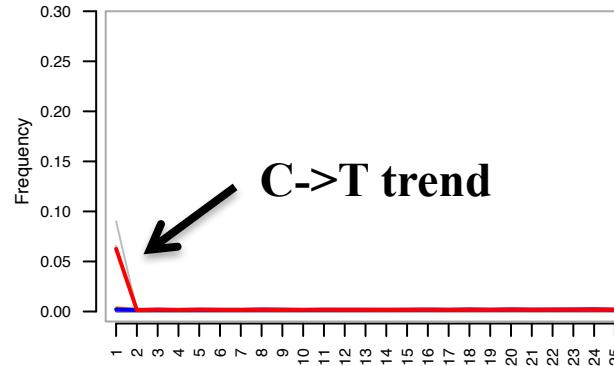


MapDamage on Controls

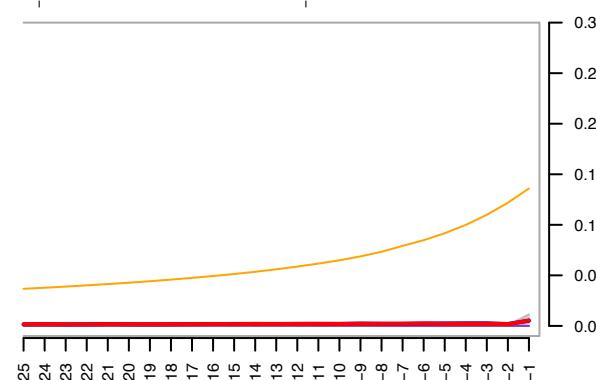
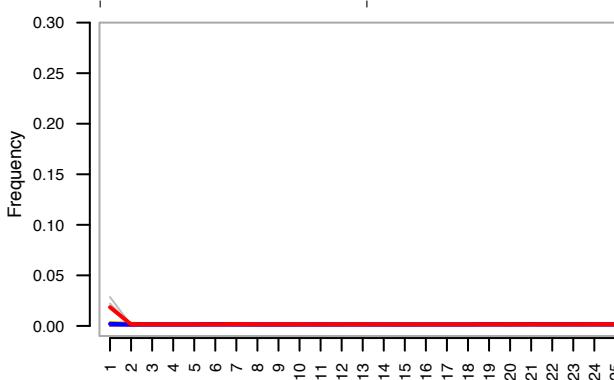


MapDamage on Moderns

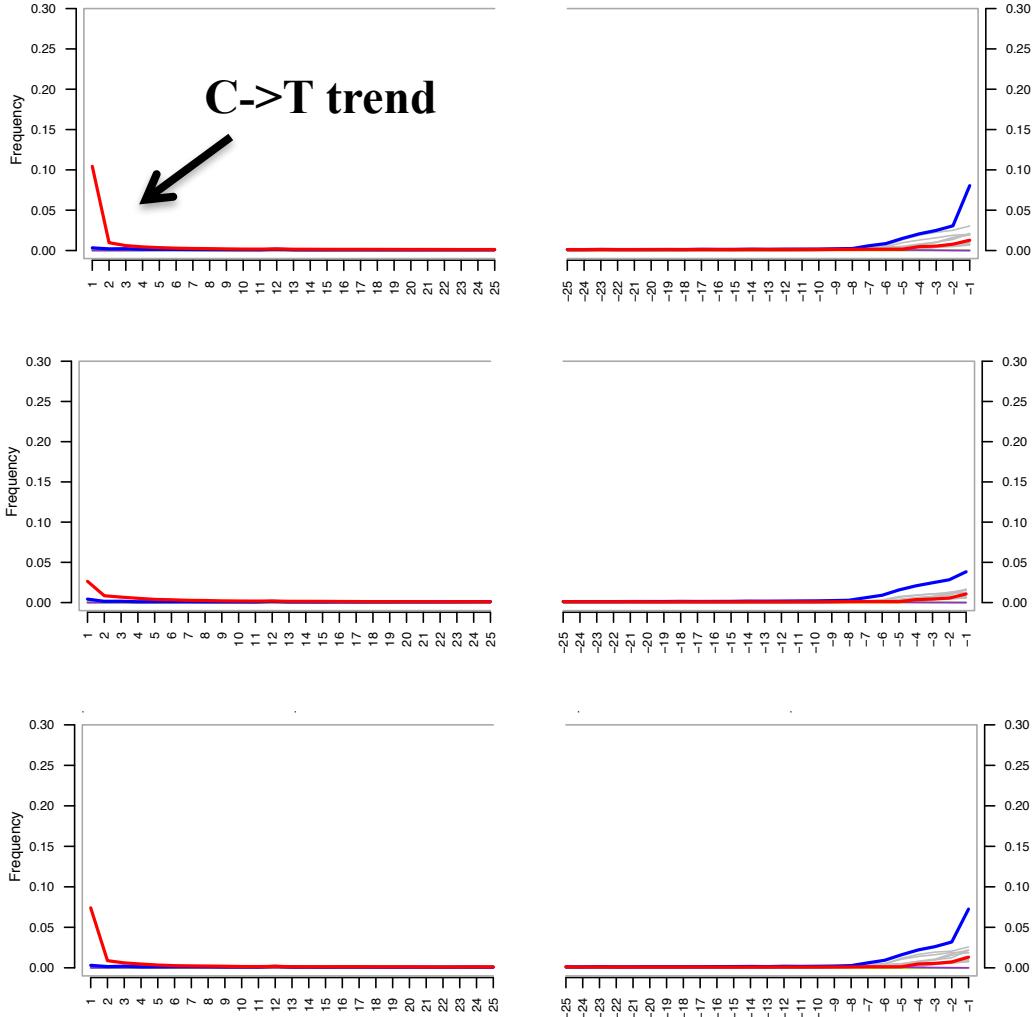
HG00097



HG00099



MapDamage on Ancient Sards. (different data-set)



Thoughts...

It's not clear what is an “ancient” C->T damage profile and what is an “modern” profile?

- Difficult to perform comparative analysis
 - MapDamage is an individual-by-individual analysis
 - Hard to look at all plots simultaneously
- Are we missing other patterns?

High-level overview of archaic

- **Step 1:** Gets “mutational patterns” from BAM files, which are flexible
- **Step 2:** Cluster and visualize samples based on “mutational patterns”.

Step 1: Get mutational patterns from BAM files

- Mutational pattern: mutation, flanking bases (we use 2), and position from the end of the read.
 - Loop through bam file and record the number of all such patterns
- Example output for Sample_1
 - AA(T->A)AA, 0, 50, 7
 - AA(T->A)AA, 10, 40, 3
 - GC(C->T)TT, 2, 48, 6

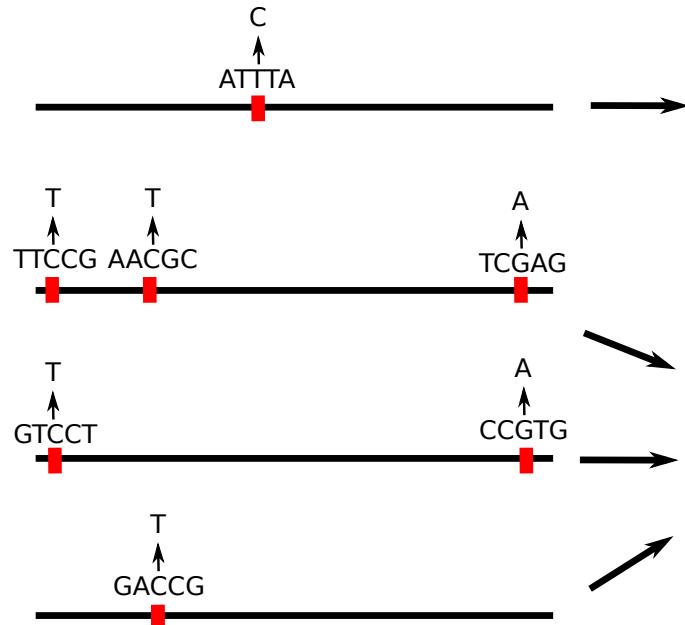
Step 2: Clustering and visualization

- STRUCTURE: if $K=2$, a African-American individual will be mixture of African ancestry and European” ancestry where ancestries are defined by a probability vector on allele frequencies
- aRchaic: if $K=2$, a contaminated individual will be mixture of ancient ancestry and modern ancestry where ancestries are defined by a probability vector on mutational patterns

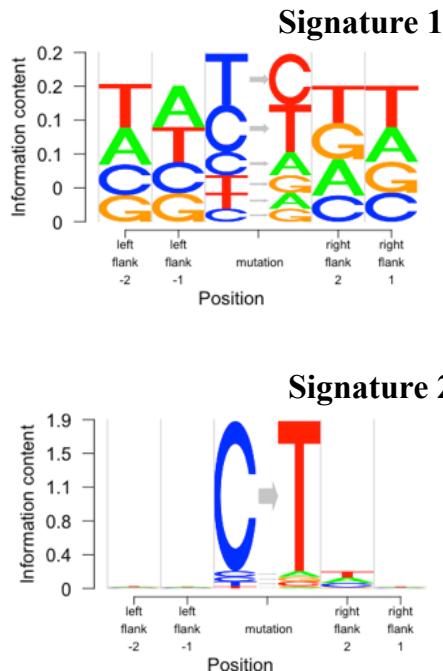


Graphical overview of archaic

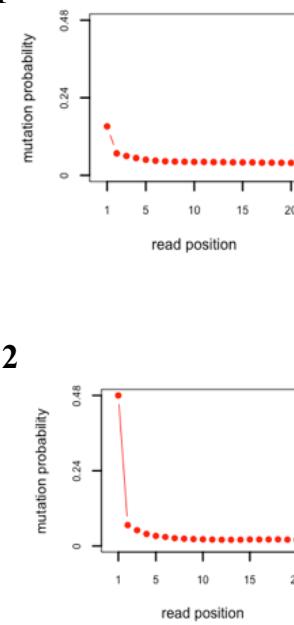
actual observed
mutational patterns



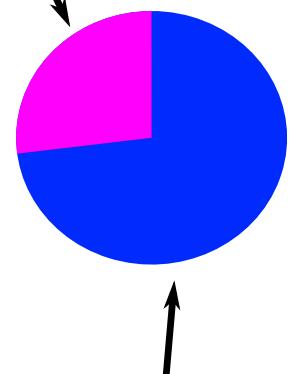
operative mutation
signatures



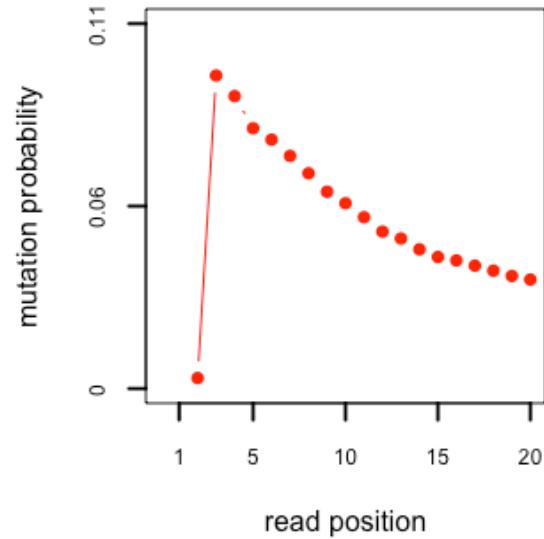
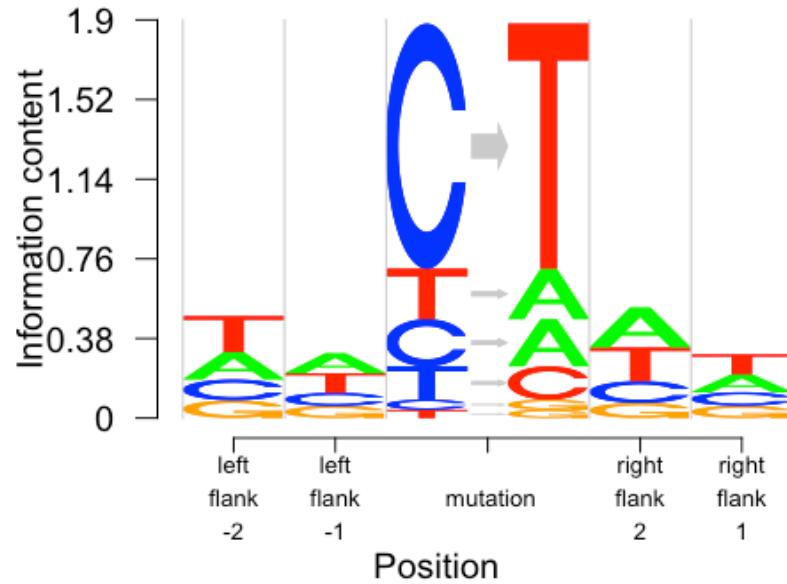
individuals are
mixtures of signatures



Signature 1



Signature 2

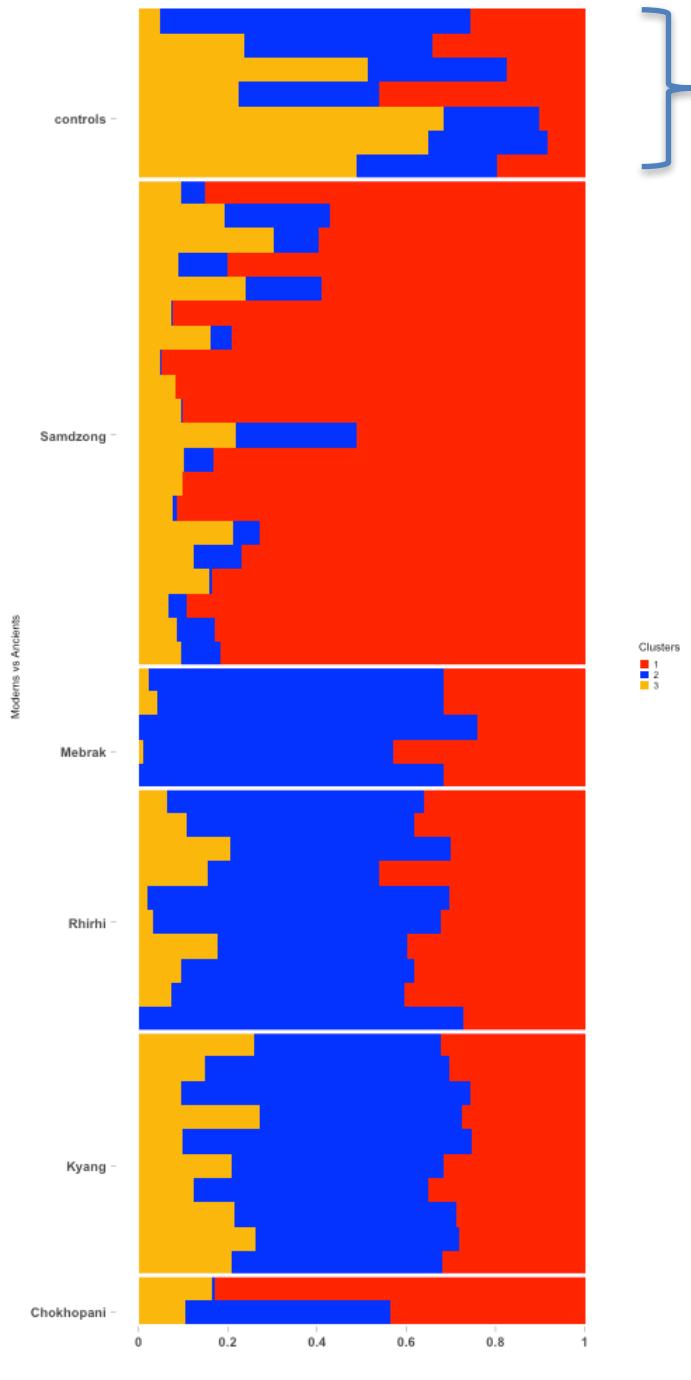


	Left flank2	Left flank 1	Right flank 1	Right flank 2
A	0.27	0.31	0.36	0.28
C	0.20	0.20	0.19	0.21
G	0.17	0.18	0.13	0.19
T	0.35	0.30	0.30	0.31

mutations	prob
C -> T	0.62
C -> A	0.12
C -> G	0.03
T -> A	0.13
T -> G	0.08
T -> C	0.02

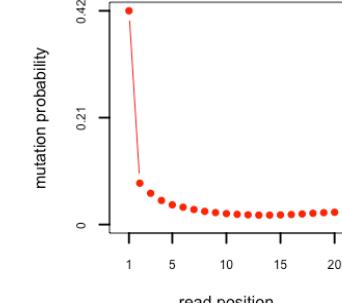
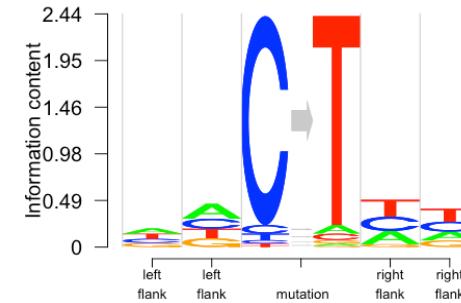
Applications of on real data

47 ancients + 7 controls from Di Rienzo lab

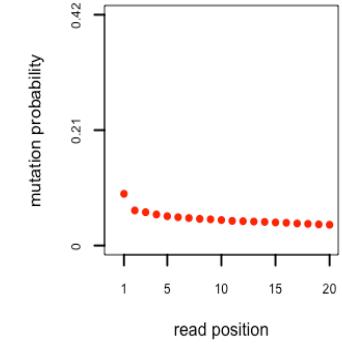
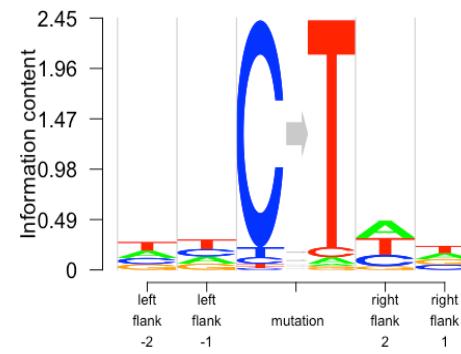


PCRNeg
Libneg3
Libneg2
Libneg1
EXN3
EXN2
EXN1

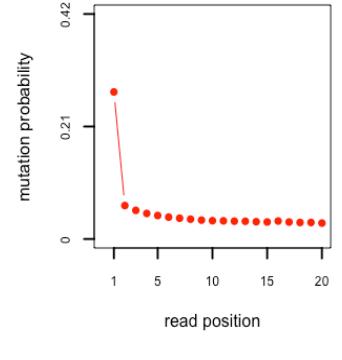
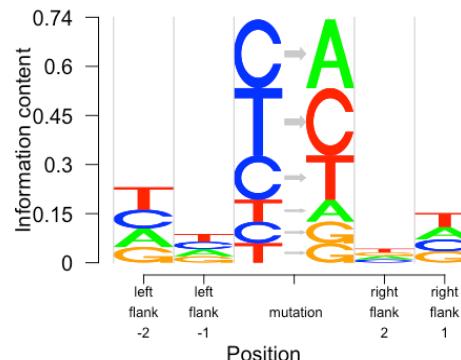
Cluster 1



Cluster 2

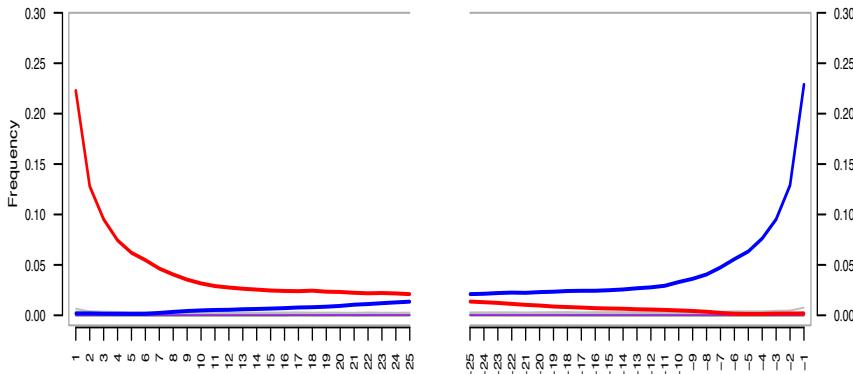


Cluster 3

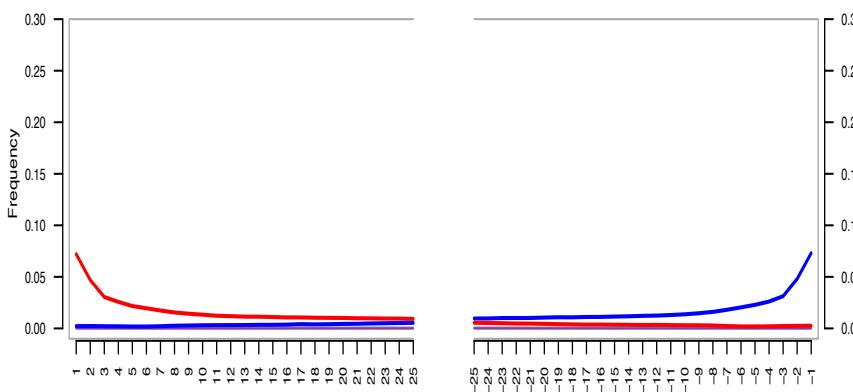


MapDamage on Ancients

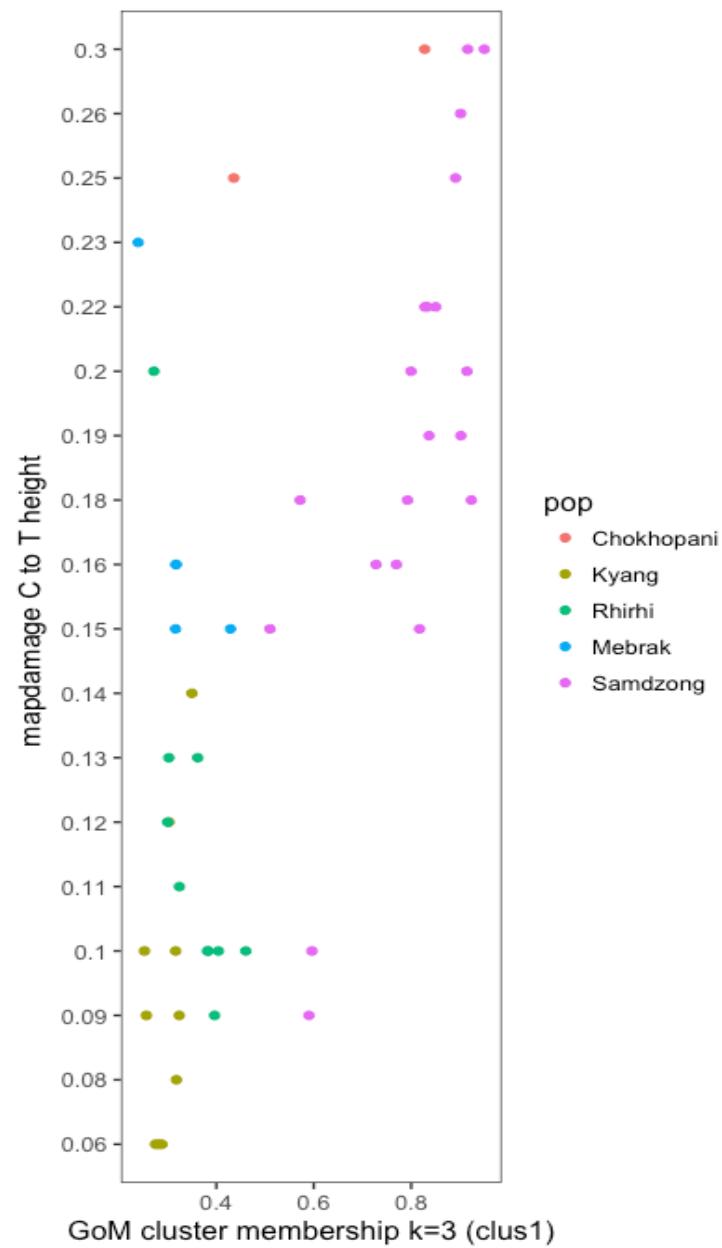
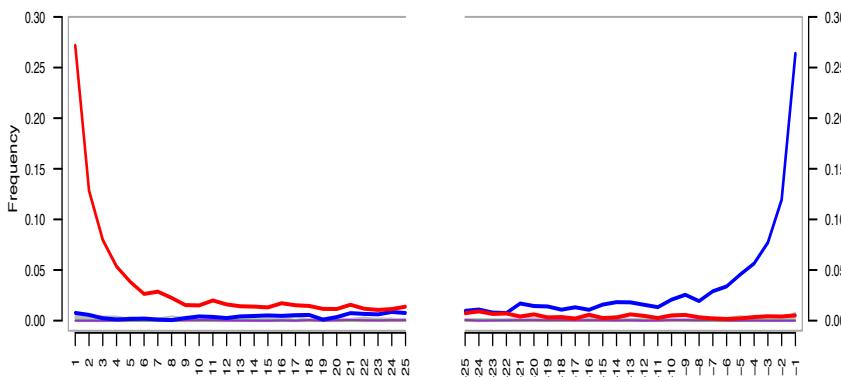
CNE1



KS20

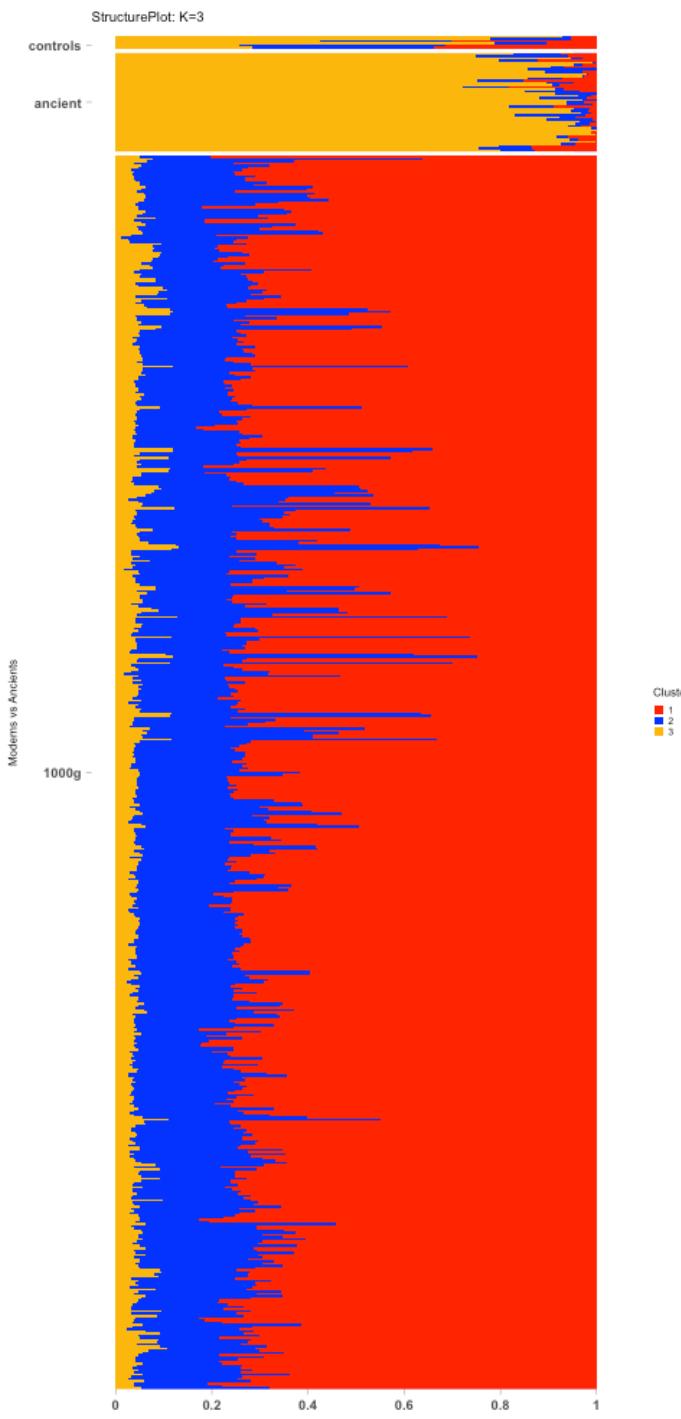


S30

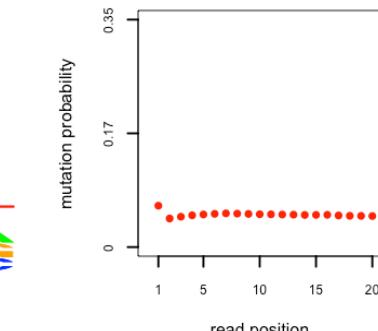
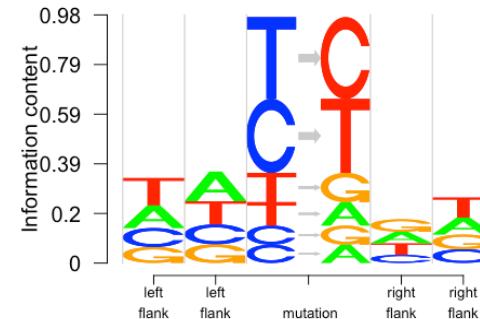


Applications of archaic on real data

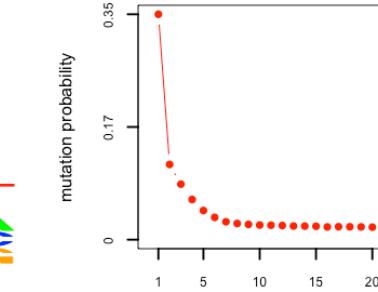
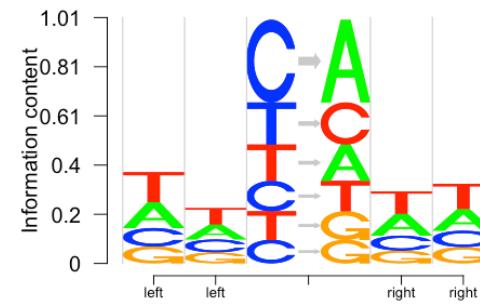
*47 ancients + 7 controls from Di Rienzo lab +
500 randomly selected 1000 Genome
modern samples of European and American
ancestry*



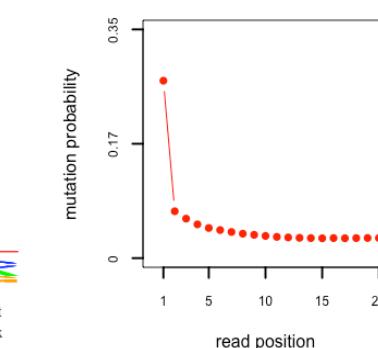
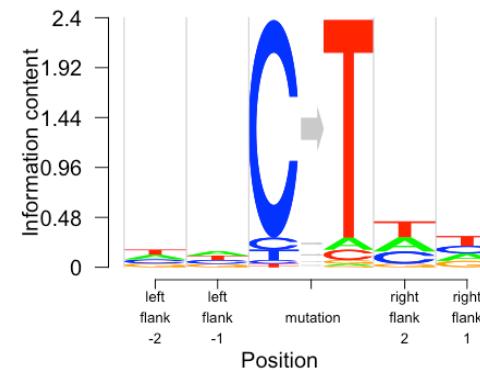
Cluster 1



Cluster 2



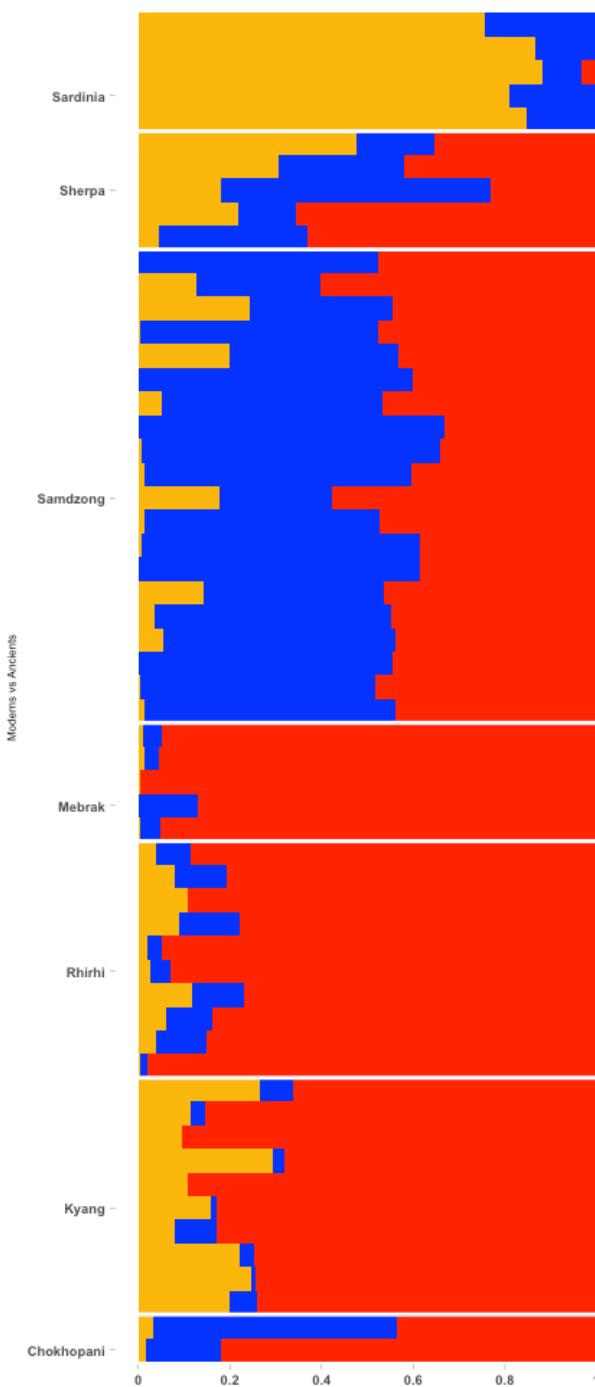
Cluster 3



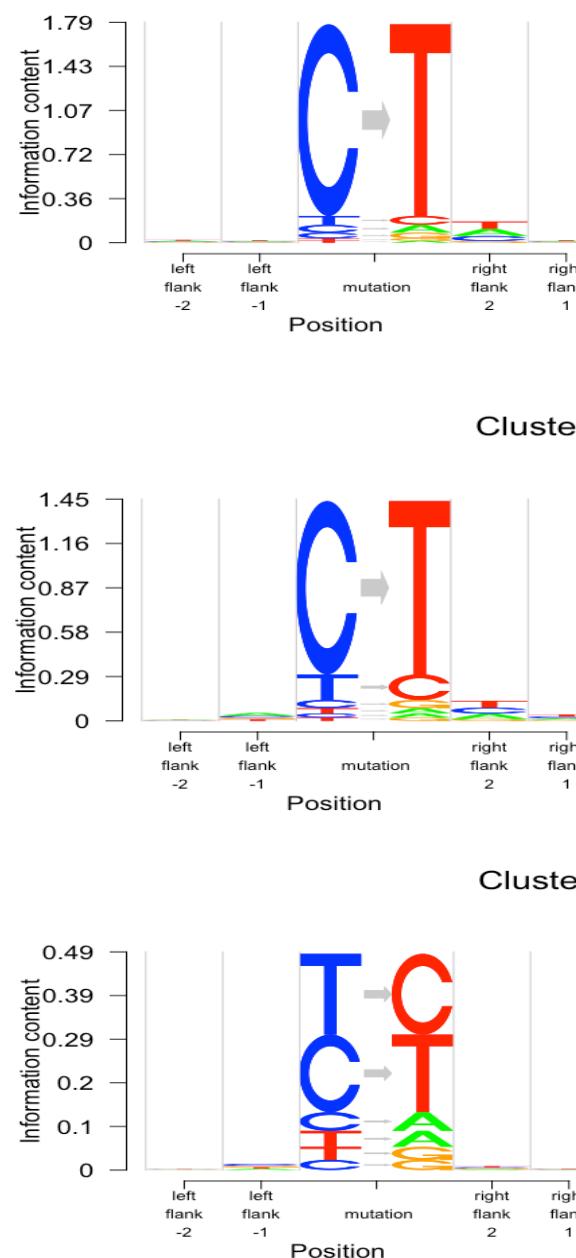
Applications of arChaic on real data

*47 ancients + 5 Sherpa samples + 5 ancient
Sardinia samples recently collected*

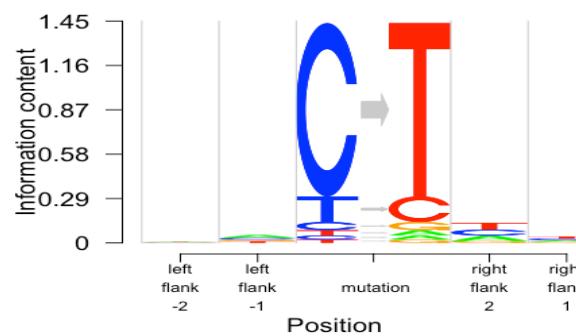
StructurePlot: K=3



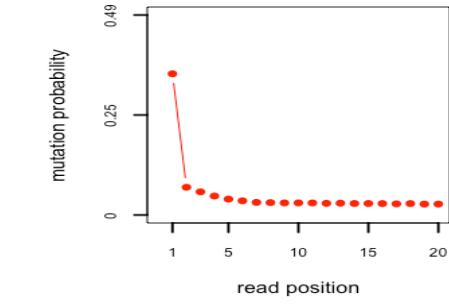
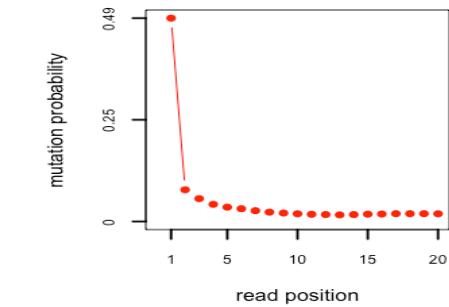
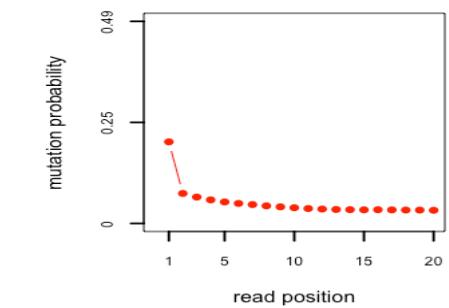
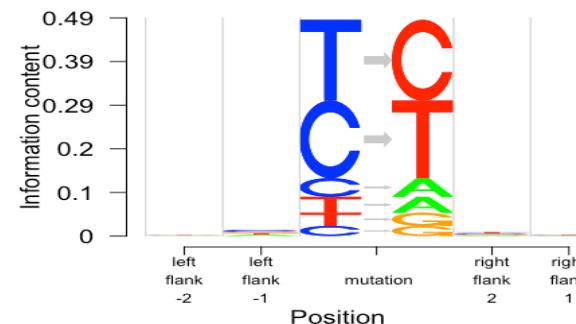
Cluster 1



Cluster 2

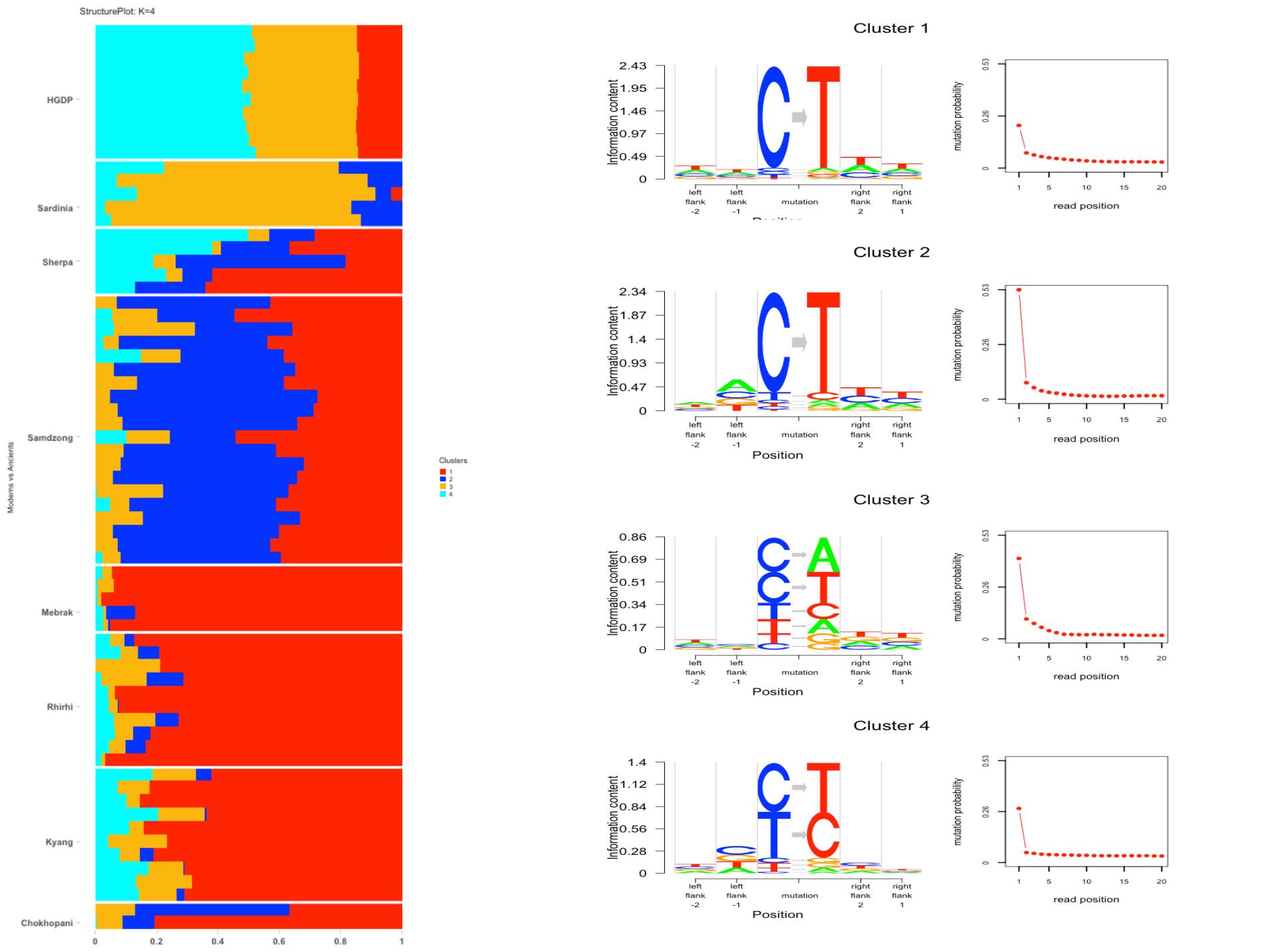


Cluster 3



Applications of archaic on real data

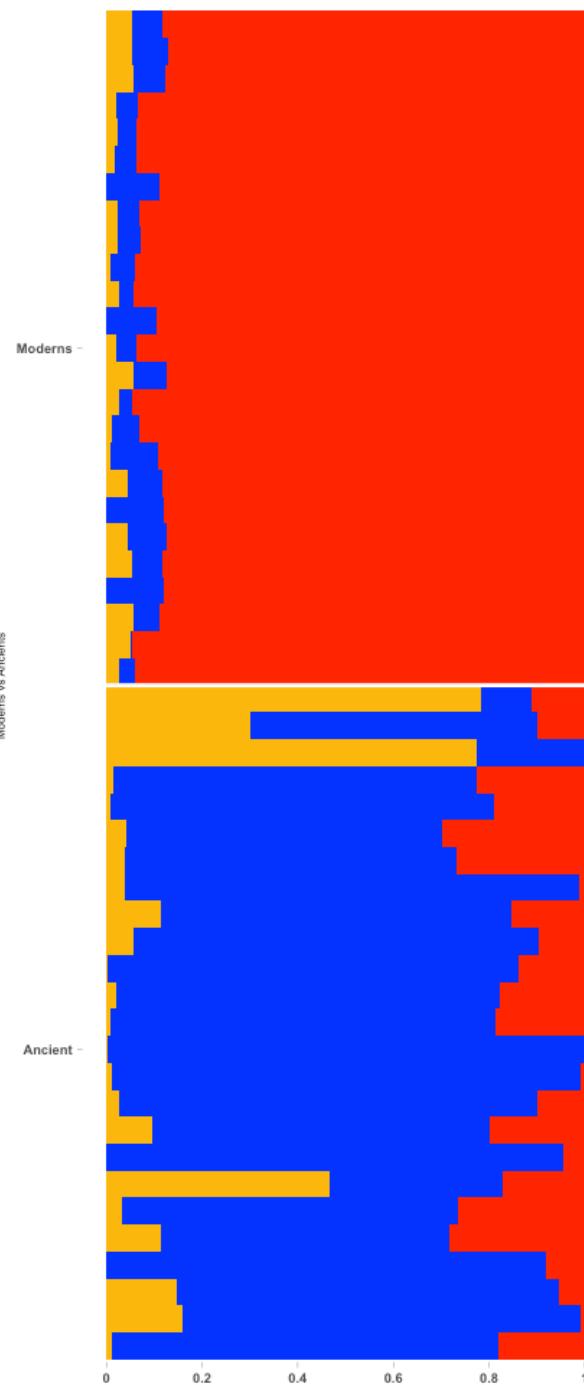
*47 ancients + 5 Sherpa samples + 5 ancient
Sardinia samples recently collected + 10
HGDP modern samples*



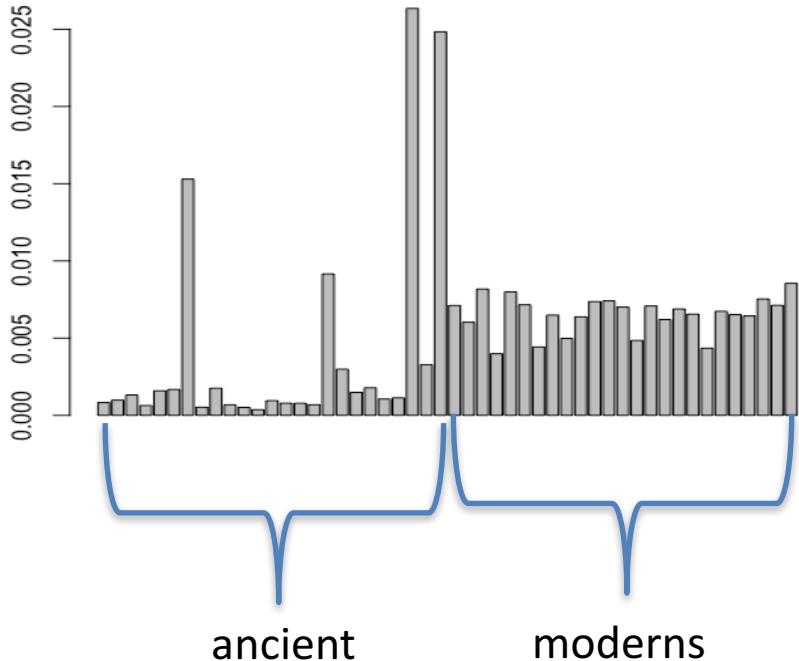
Applications of archaic on real data

*25 ancients + 25 moderns from the ancient
data collected by John Lindo (DiRienzo Lab)*

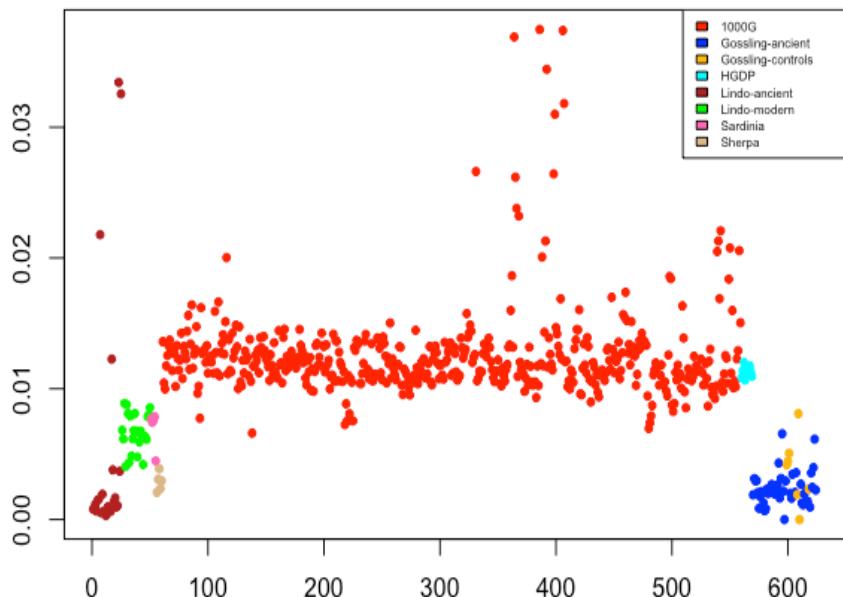
StructurePlot: K=3



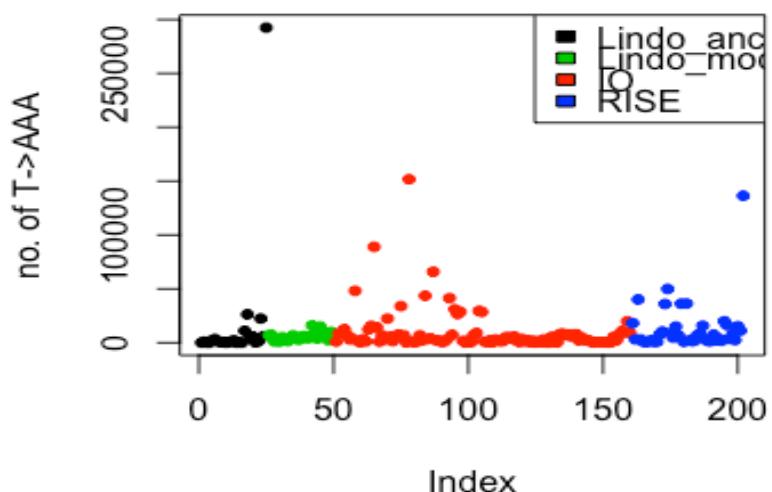
Proportion of T->AAA (each sample)



prop. of T->AAA across aDNA sources



Number of T->AAA (each sample)



Applications of archaic on real data

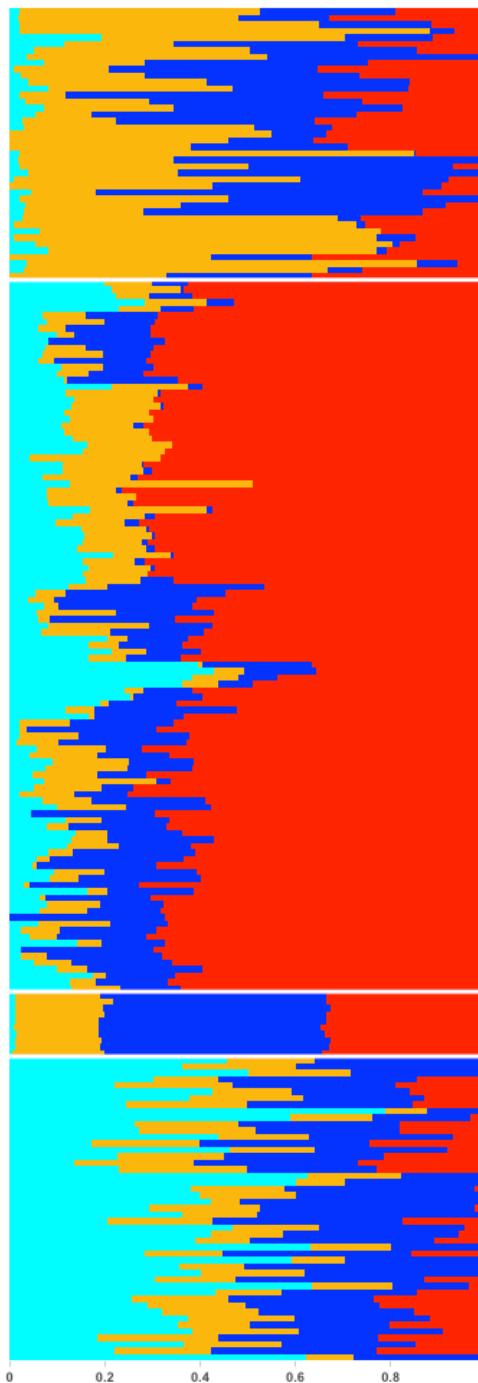


*Allentoft + Haak et al published ancient data
+ Di Rienzo data + 10 HGDP modern samples*

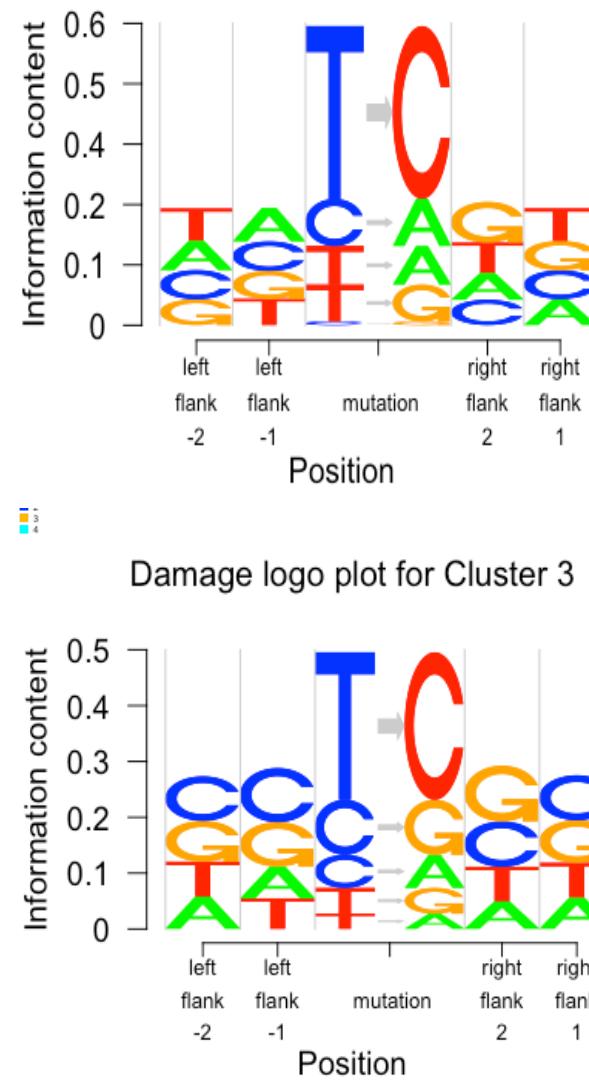
*We show the lab specific effects on
mutational patterns.*

*We remove the C-> T patterns and the read
position information and ran the model on
remaining mutational patterns.*

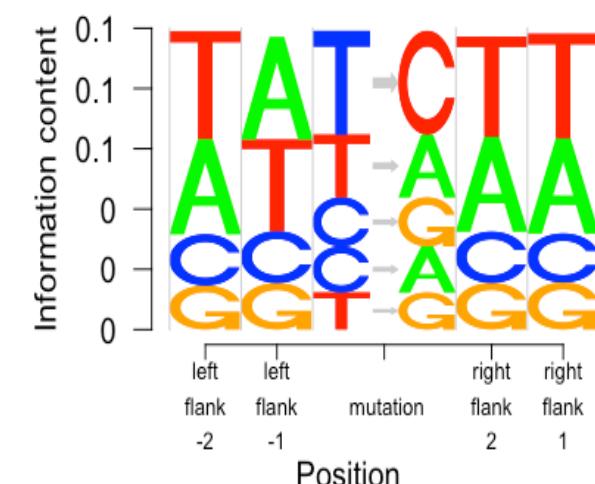
RISE



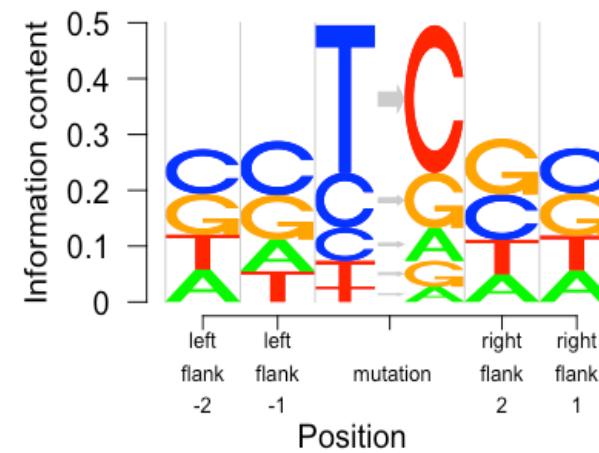
Damage logo plot for Cluster 1



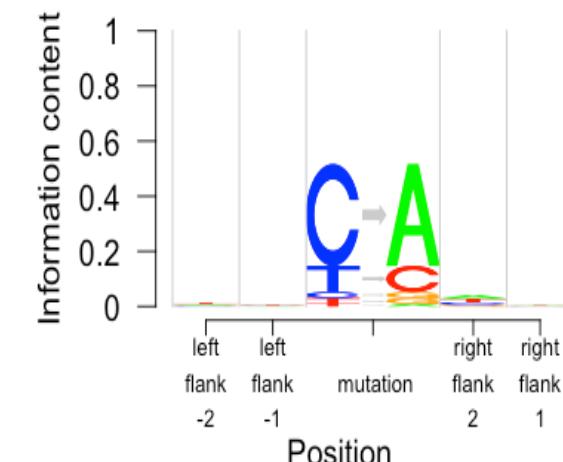
Damage logo plot for Cluster 2



Damage logo plot for Cluster 3



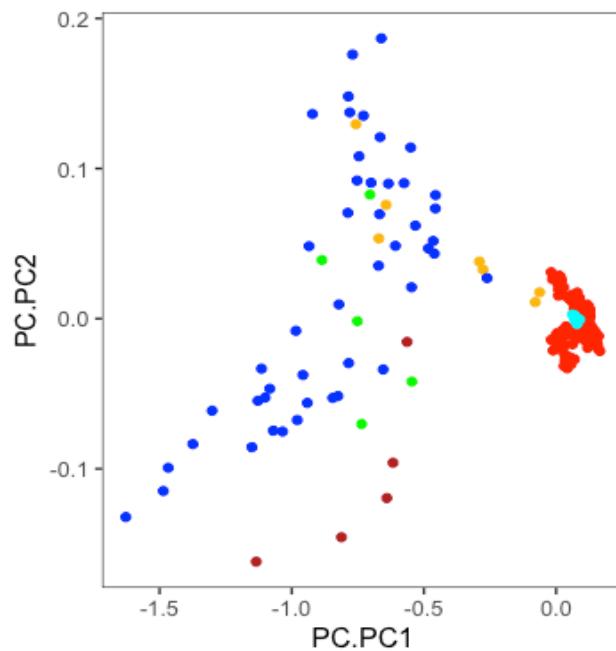
Damage logo plot for Cluster 4



Applications of archaic on real data

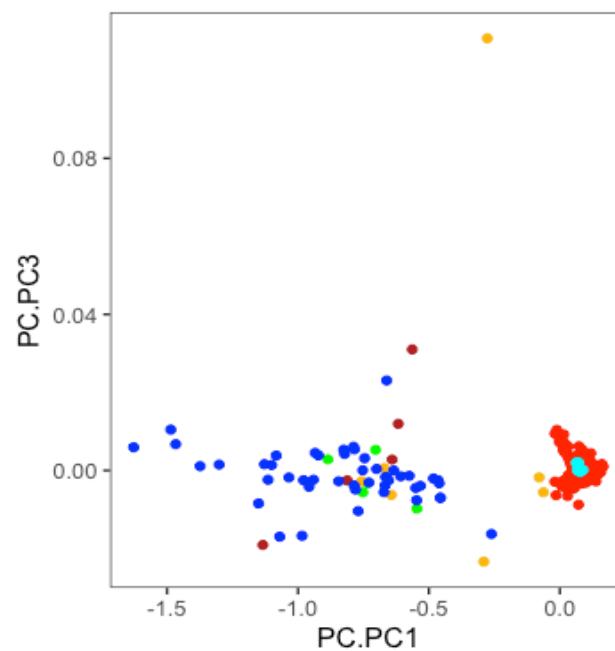
Sherpa + Gossling + HGDP + 100 randomly selected samples from 1000 Genomes + ancient Sardinian data

We focus on only known damage signatures – C->T along with the position of the mutation from end of the read (from position 0 to position 10)



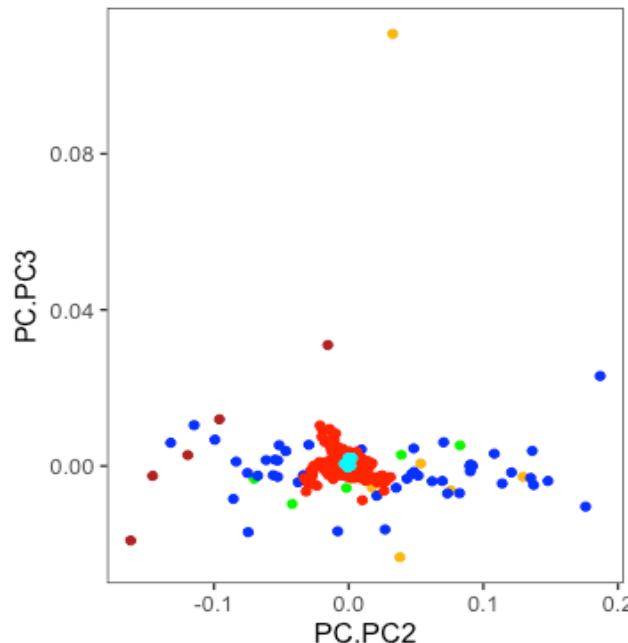
factor(labels)

- 1000G
- gosling-ancient
- gosling-control
- hgdp
- sardinia
- sherpa



factor(labels)

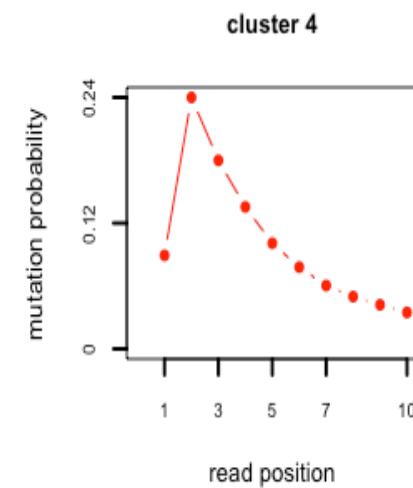
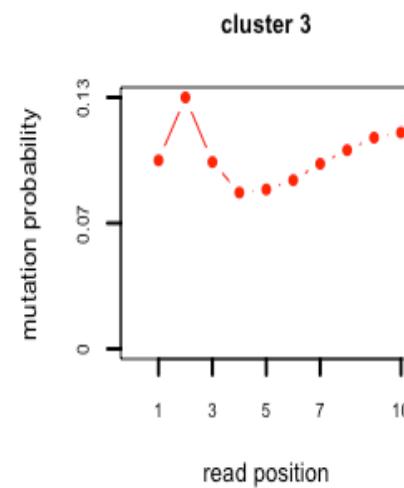
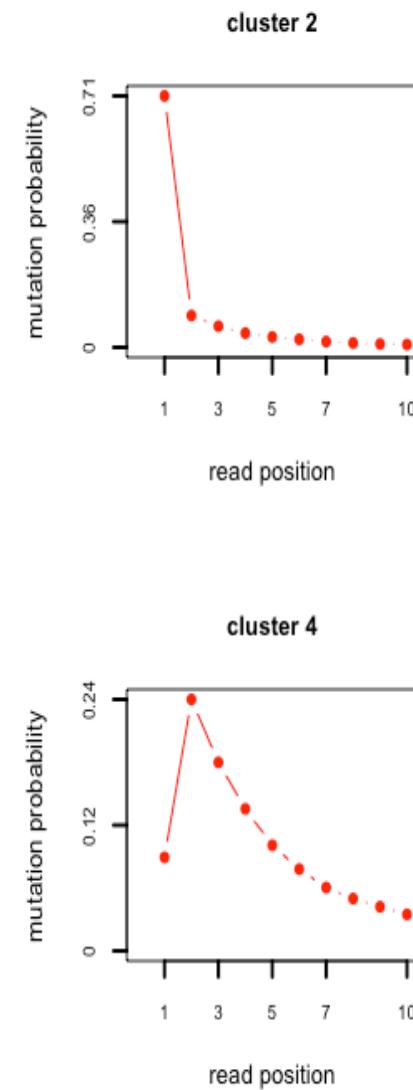
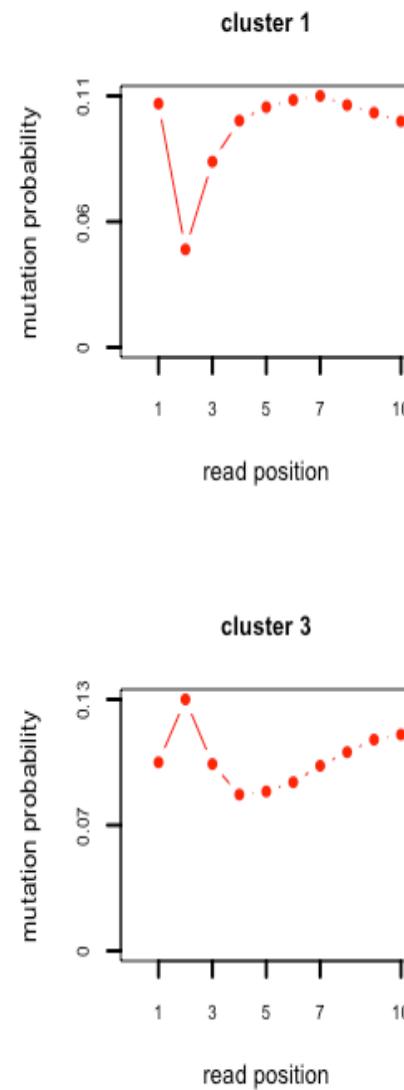
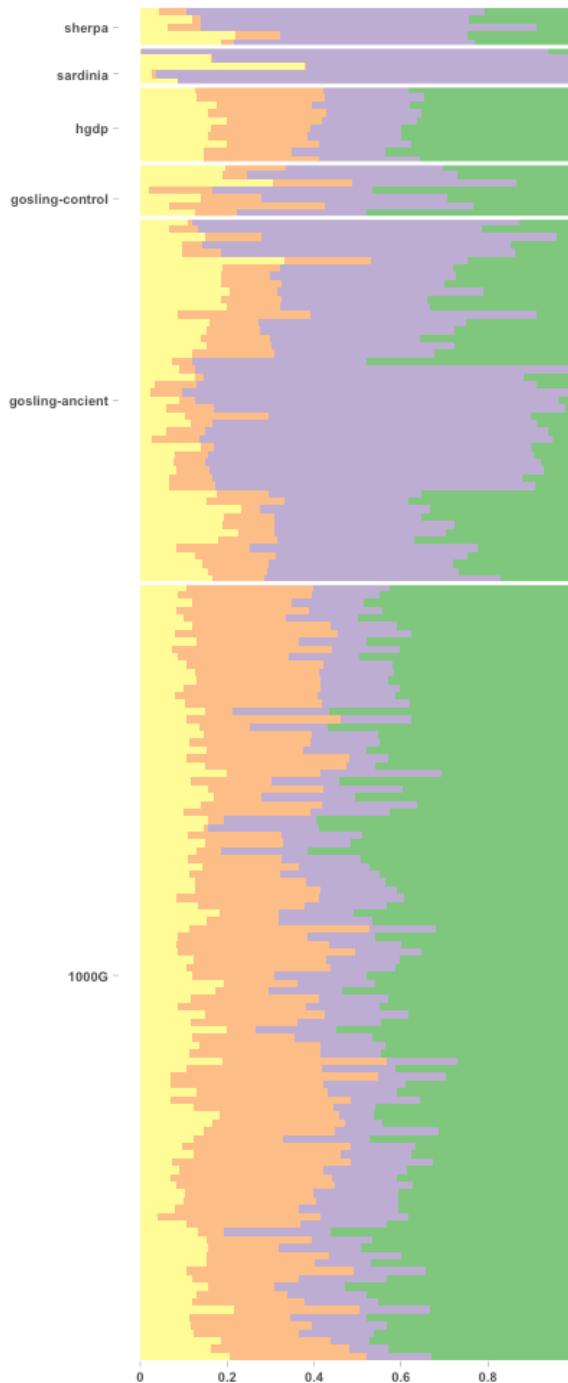
- 1000G
- gosling-ancient
- gosling-control
- hgdp
- sardinia
- sherpa



factor(labels)

- 1000G
- gosling-ancient
- gosling-control
- hgdp
- sardinia
- sherpa

Tissue type

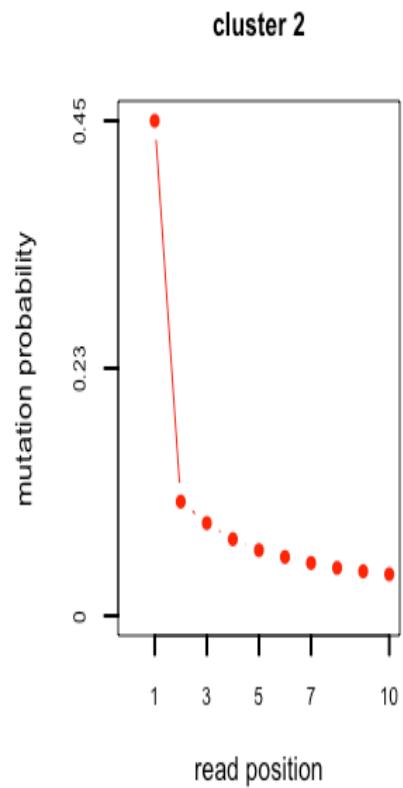
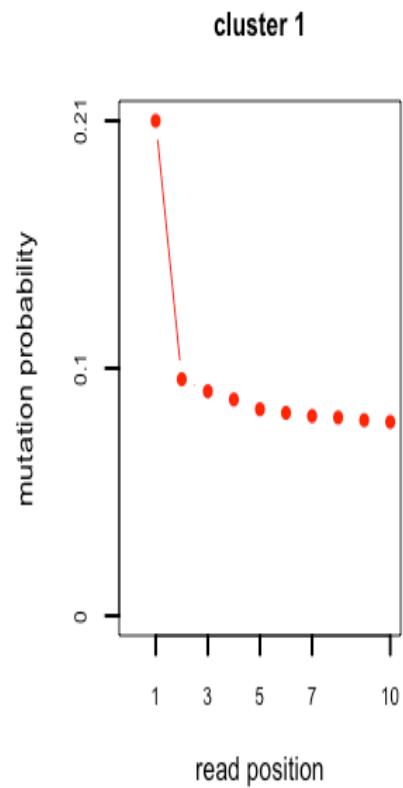
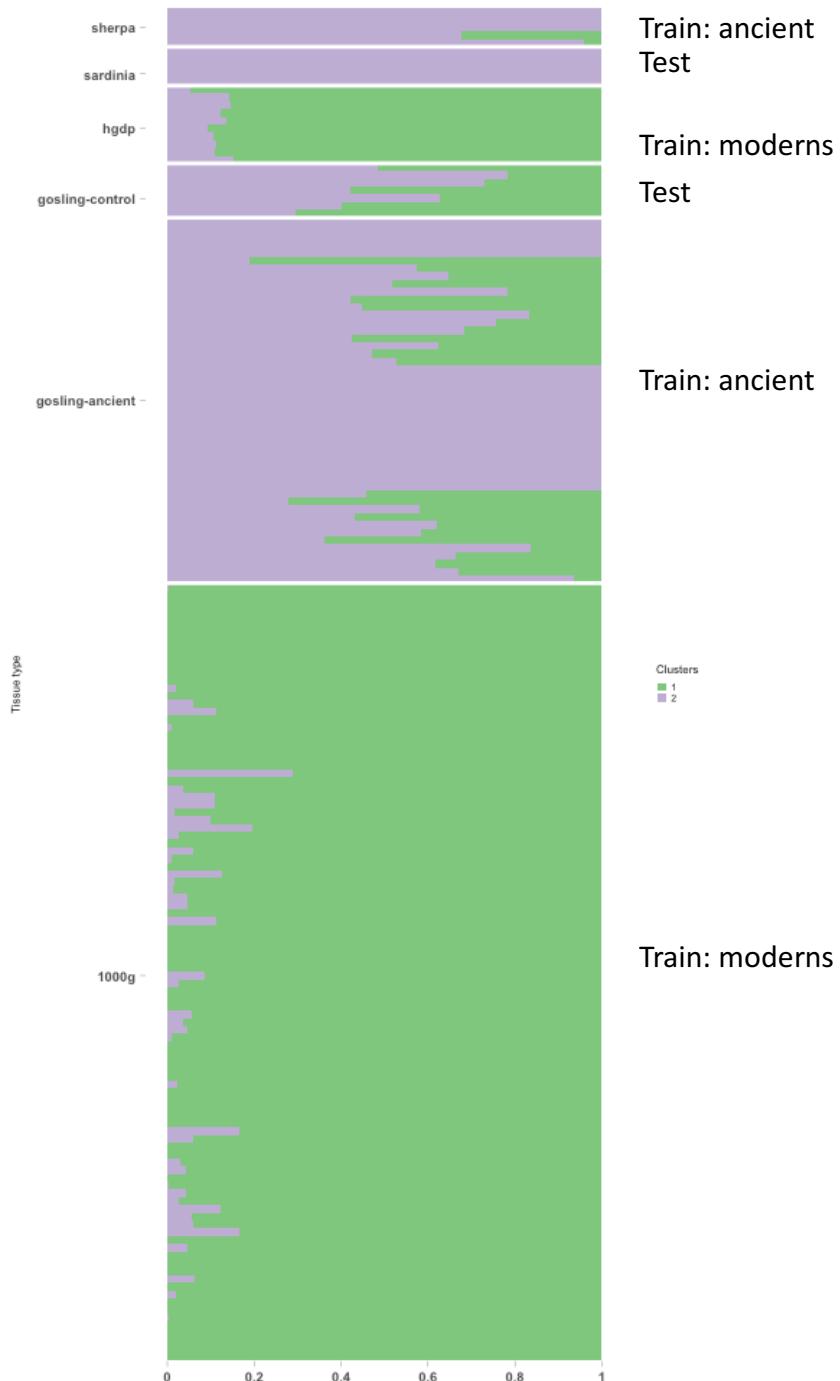


Applications of **aRchaic**
In classifying between
moderns and ancient
samples

Consensus ancient data panel : Sherpa + the Gossling ancient.

Consensus modern data panel: 1000 Genomes + HGDP moderns.

- We train the moderns and the ancient panels using the above data and test on the Sardinian ancients and the Gossling controls data.
- We apply Support Vector machines for classifying the test samples and generate a probability for each test sample to come from consensus moderns and equivalently, from consensus ancient panels.
- For Sardinian samples (total 5), 2 samples were classified as ancients and 3 as moderns.
- For Gossling controls data (total 7), 2 were classified strongly as ancients (EXN3 and Libneg3) and the rest as moderns.



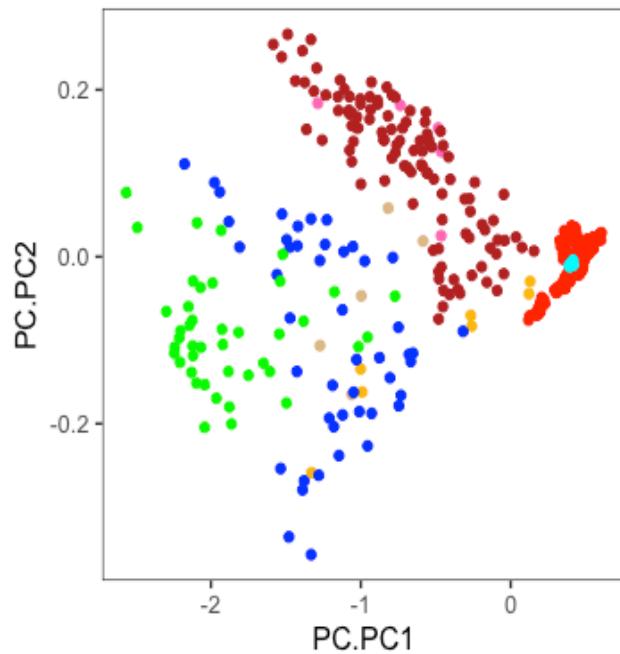
R package aRchaic

Developmental version :

<https://github.com/kkdey/aRchaic>

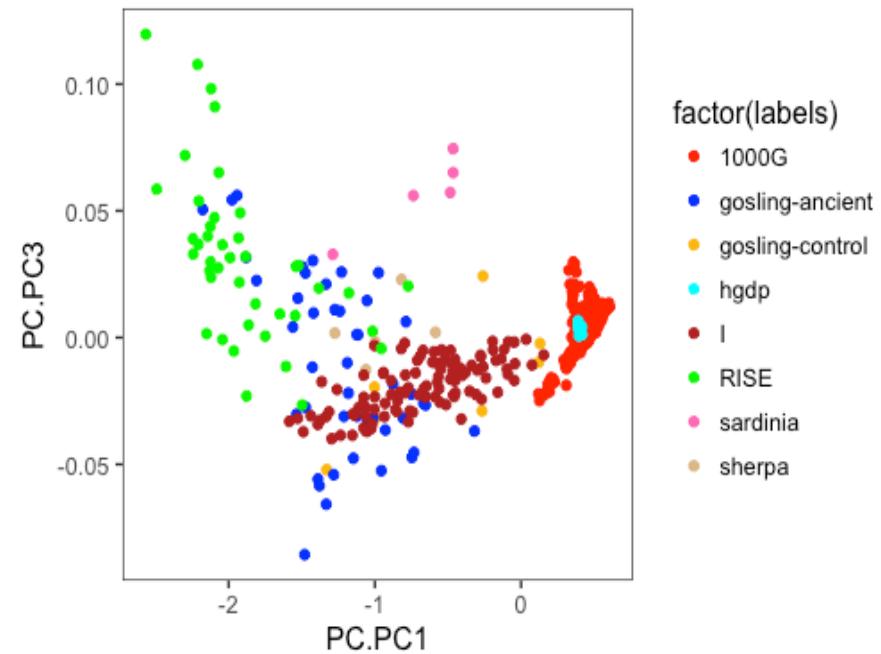
The codes and the scripts for generating
the plots and also additional plots for
other values of K can be found at

<https://github.com/halasadi/ancient-damage>



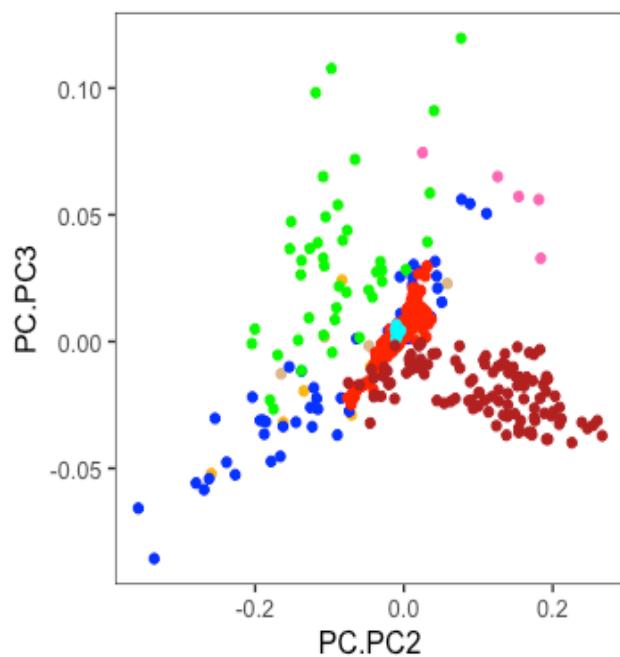
factor(labels)

- 1000G
- gosling-ancient
- gosling-control
- hgdp
- I
- RISE
- sardinia
- sherpa



factor(labels)

- 1000G
- gosling-ancient
- gosling-control
- hgdp
- I
- RISE
- sardinia
- sherpa



factor(labels)

- 1000G
- gosling-ancient
- gosling-control
- hgdp
- I
- RISE
- sardinia
- sherpa