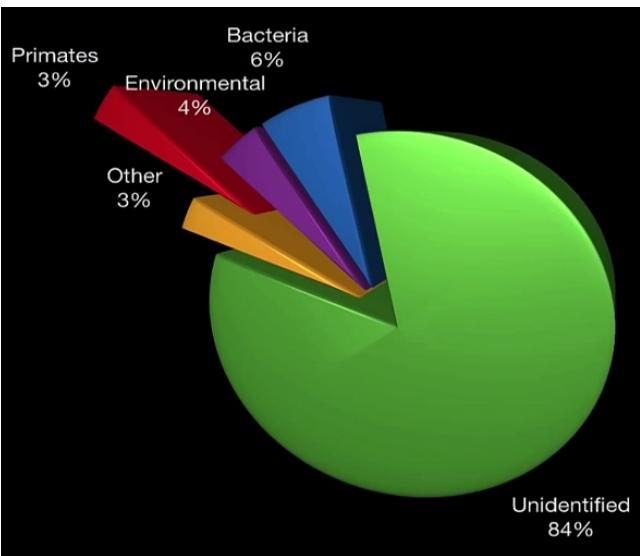


aRchaic

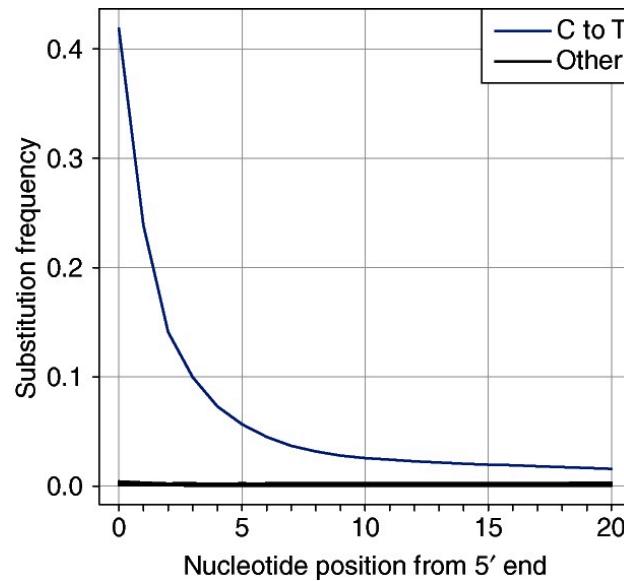
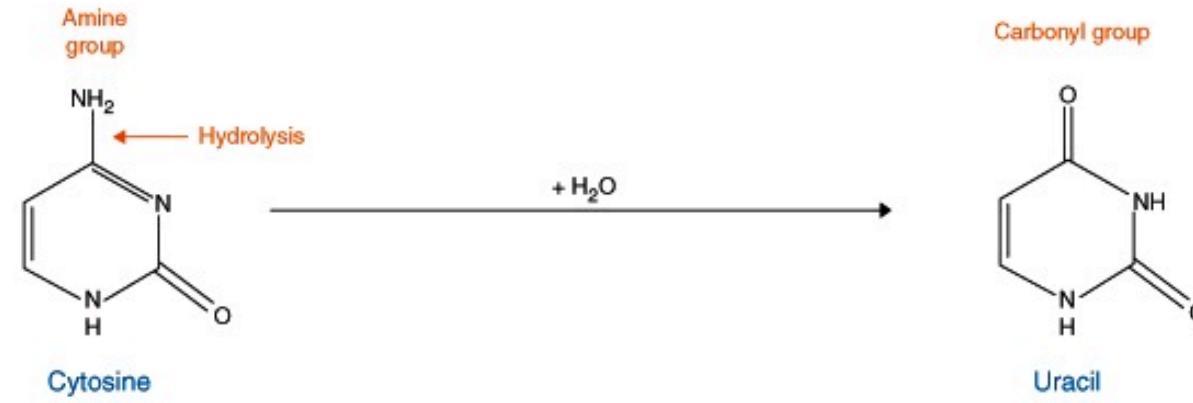


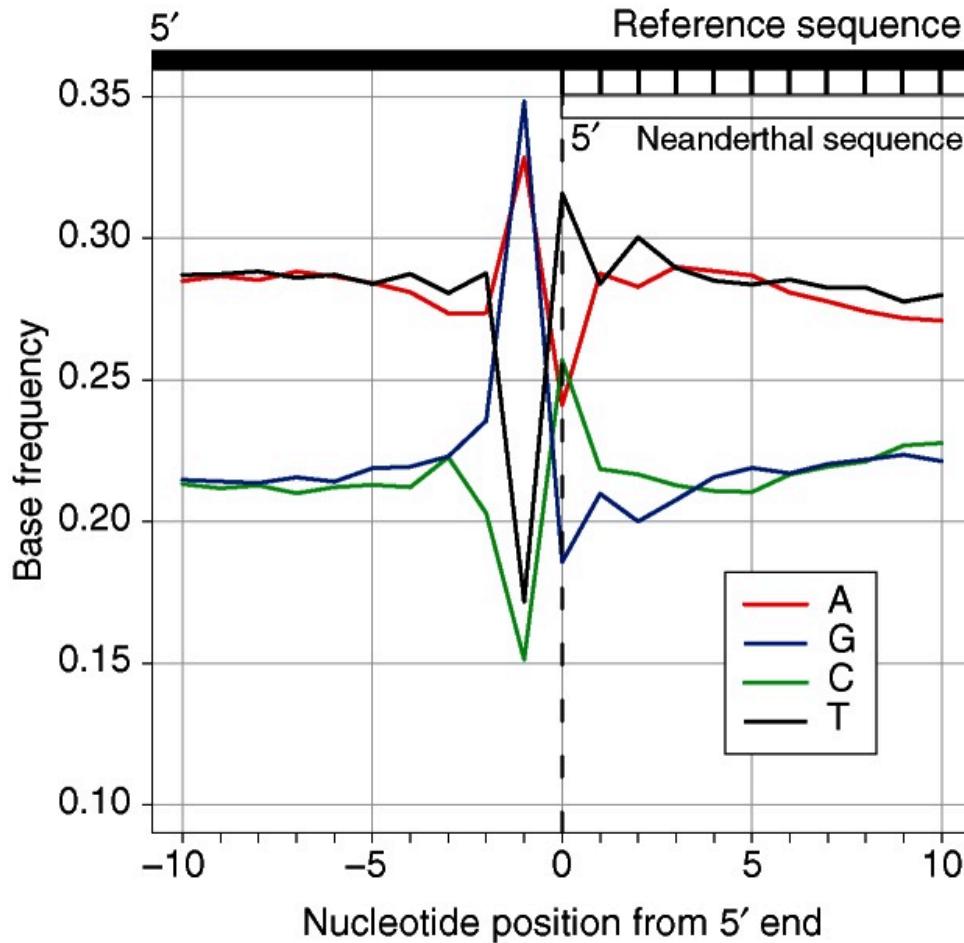
An R package for clustering and
visualizing ancient dna signatures

By Hussein Al-Asadi & Kushal K Dey



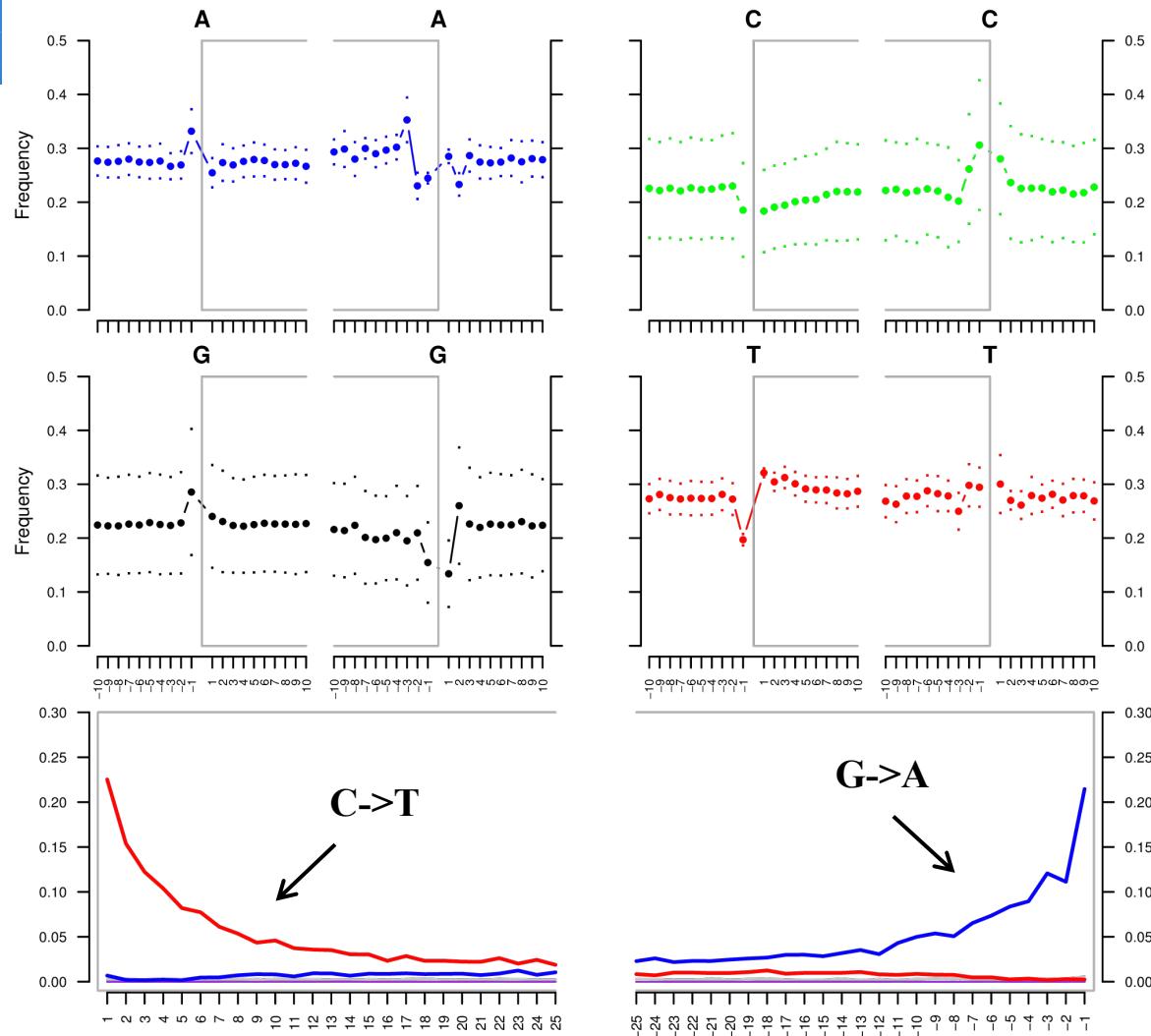
- DNA degradation and damage
- Bacterial DNA
- Human contamination





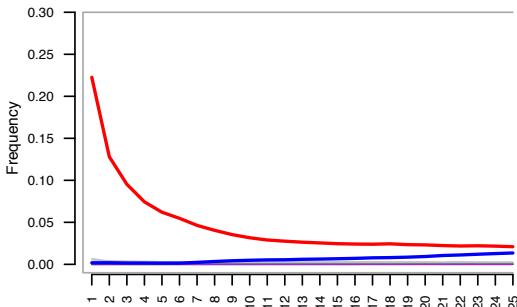
mapDamage

mapDamage: tracking and quantifying
damage patterns in ancient DNA sequences

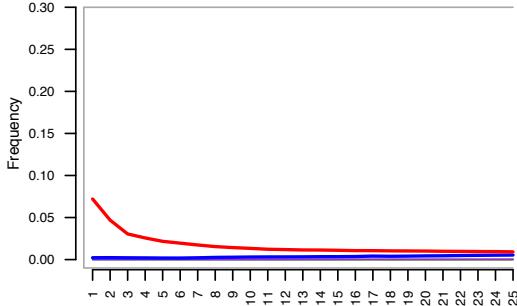


MapDamage on Ancients

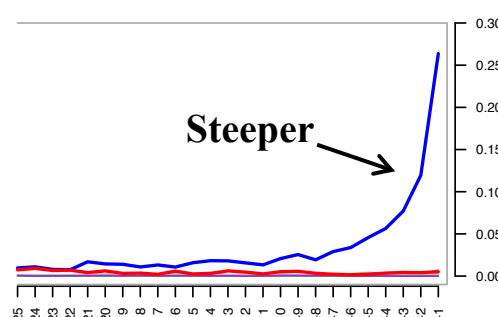
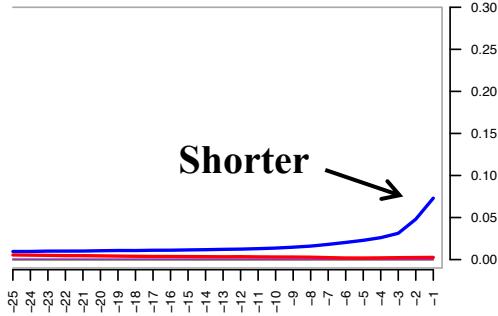
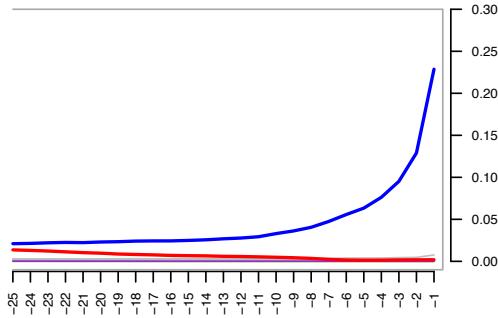
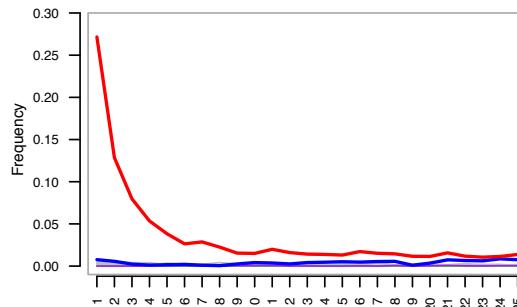
CNE1



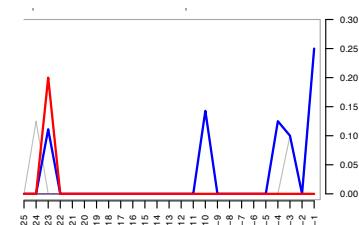
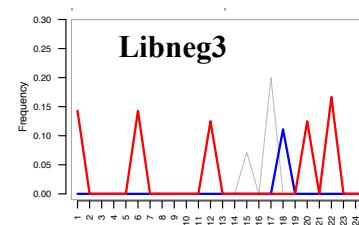
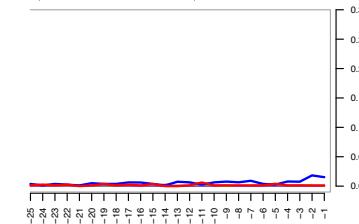
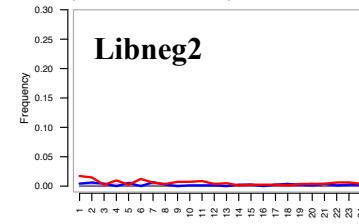
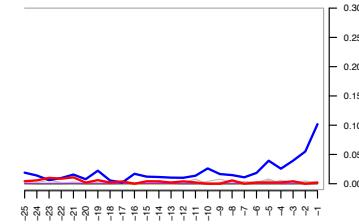
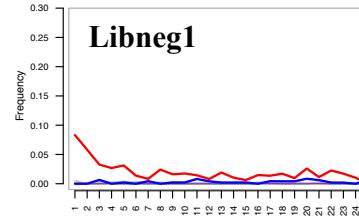
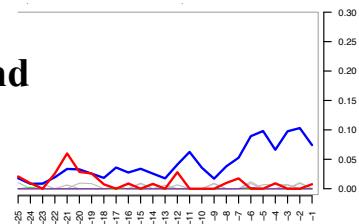
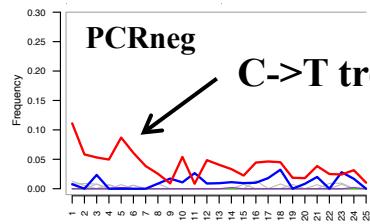
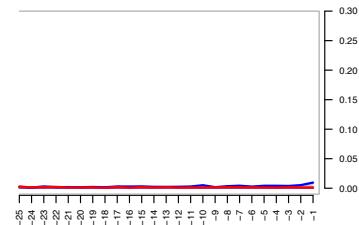
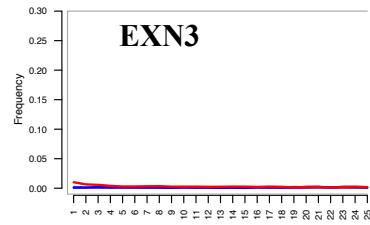
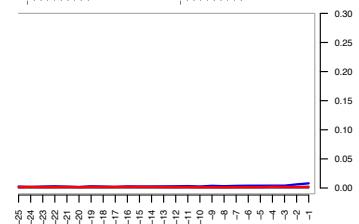
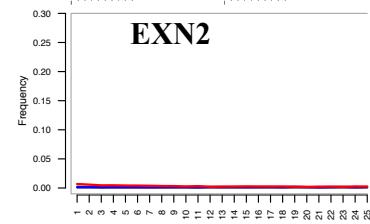
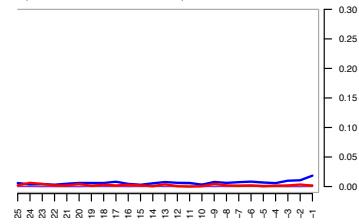
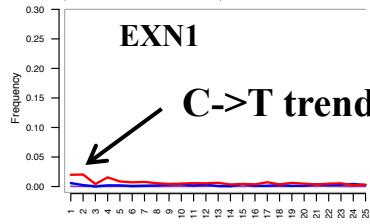
KS20



S30

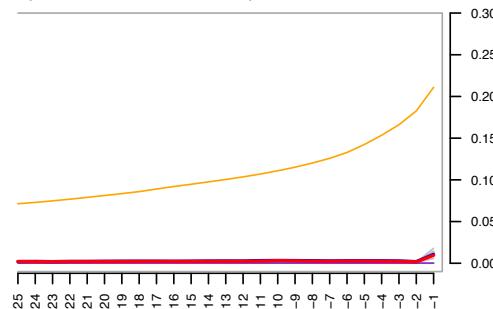
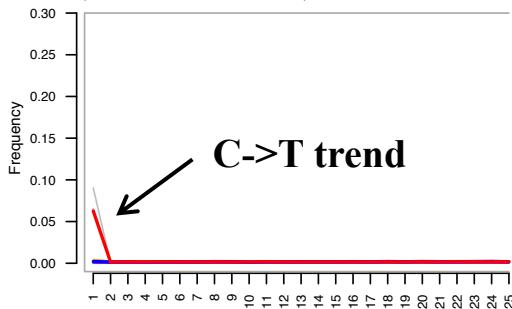


Control

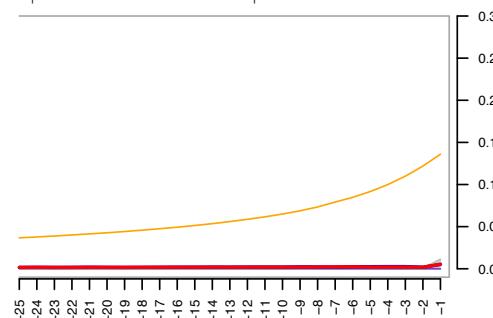
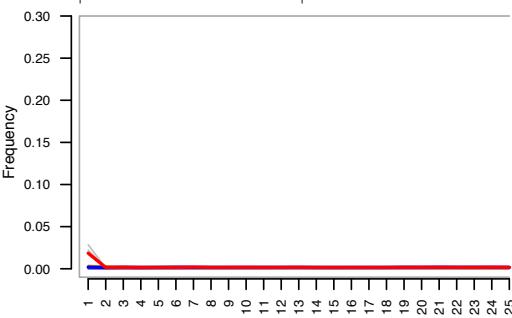


MapDamage on Moderns

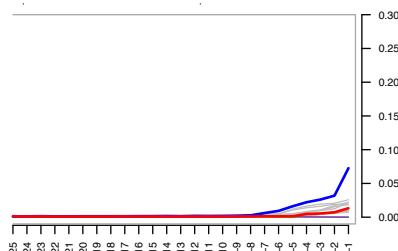
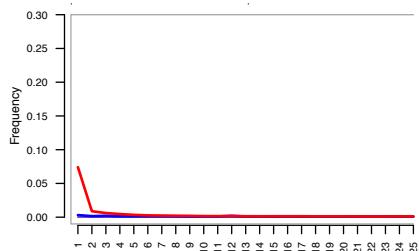
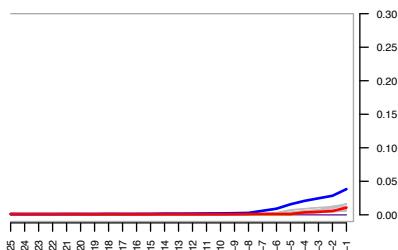
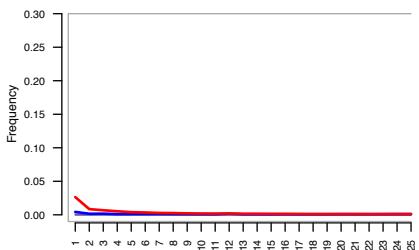
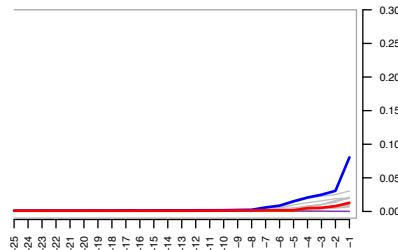
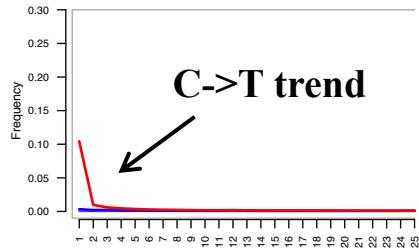
HG00097



HG00099



MapDamage on Ancient Sards. (different data-set)



Thoughts...

- It's not clear what is a C->T damage profile is and what is a damage profile typical of a modern individual is.
- Difficult to perform joint analysis
 - MapDamage is an individual-by-individual analysis
 - Hard to look at all plots simultaneously
- Are we missing other patterns?

High-level overview of

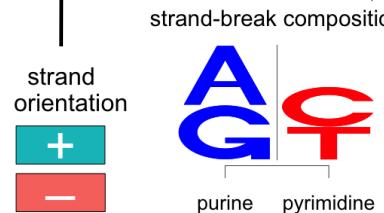
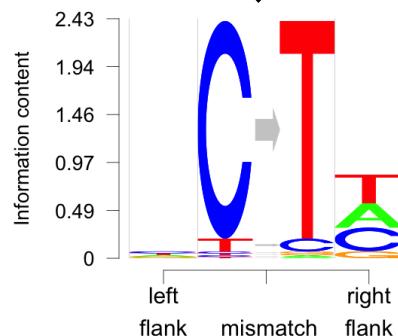
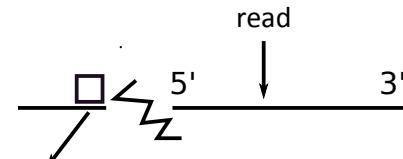
- **Step 1:** Gets “mismatch patterns” (i.e. features) from BAM files, which are flexible
- **Step 2:** Cluster and visualize samples based on “mismatch patterns”.

Step 1: Get mismatch patterns from BAM files

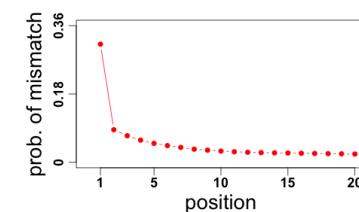
- Mutational pattern:
 - mutation
 - flanking bases (we use 1)
 - position from the end of the read
 - strand orientation
 - base pair immediately outside 5' strand break.
- Example output for Sample_1
 - A(T->C)A, 0, 50, +, A, 2
 - A(T->A)A, 10, 40, - , T, 1
 - C(C->T)T, 2, 48, +, G, 6

C → T	T → C	T → G	C → A	C → G	T → A
0.92	0.054	0.007	0.007	0.007	0.005

strand-orientation	Percentage
+	0.5
-	0.5



Base	Percentage
A	0.36
C	0.18
G	0.32
T	0.14



Base	Percentage
A	0.24
C	0.28
G	0.23
T	0.25

Base	Percentage
A	0.29
C	0.28
G	0.08
T	0.35

Probability of mismatch across the read

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 a “**damaged**” ancestry and “**un-damaged**” ancestry where ancestries are **defined** by a probability vector on mutational patterns

Just Notation

- $i = \{1, \dots, I\}$ I individuals (bam files)
- Each individual has J_i mismatch patterns.
- Let M denote the number of features, here $M = 6$.
- $x_{i,j,l}$ where $j = \{1, \dots, J_i\}$, $l = \{1, \dots, M\}$
 - e.g. $x_{i,j,1} = A$ if 1 is the mutational feature corresponding to the left flanking base.

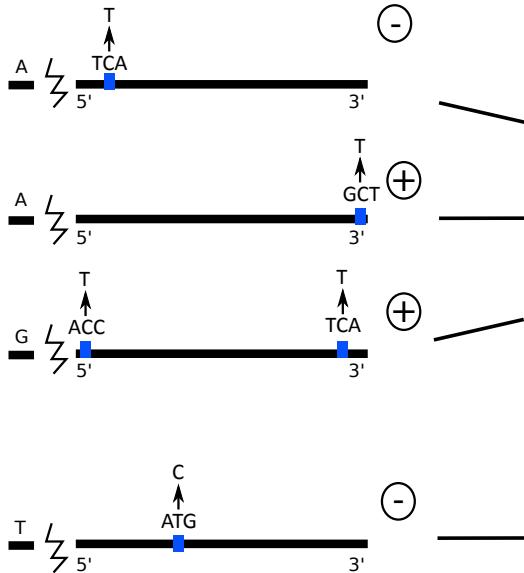
Generative Model

- (1) Generate $z_{i,j} \sim \text{Multinomial}(1, q_i)$ where $z_{i,j} \in 1, \dots, K$ and $i = 1, \dots, I$ and $j = 1, \dots, J_i$
- (2) For each $l = 1, \dots, M$, generate $x_{i,j,l} \sim \text{Multinomial}(1, f_{z_{i,j},l})$

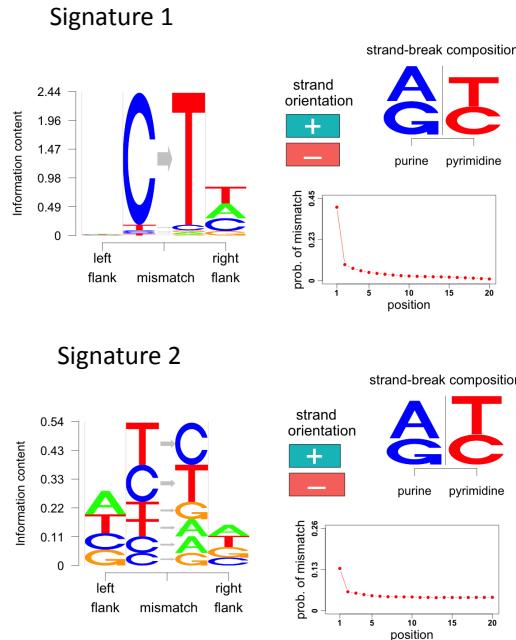
Graphical overview of

aRchaic

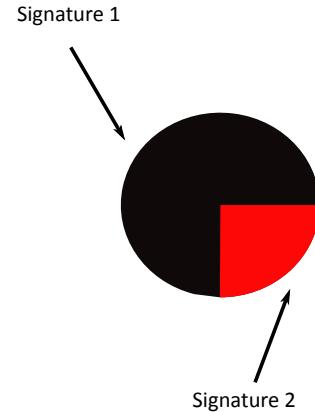
observed mismatch patterns



operative mismatch signatures



individuals are a mixture of signatures



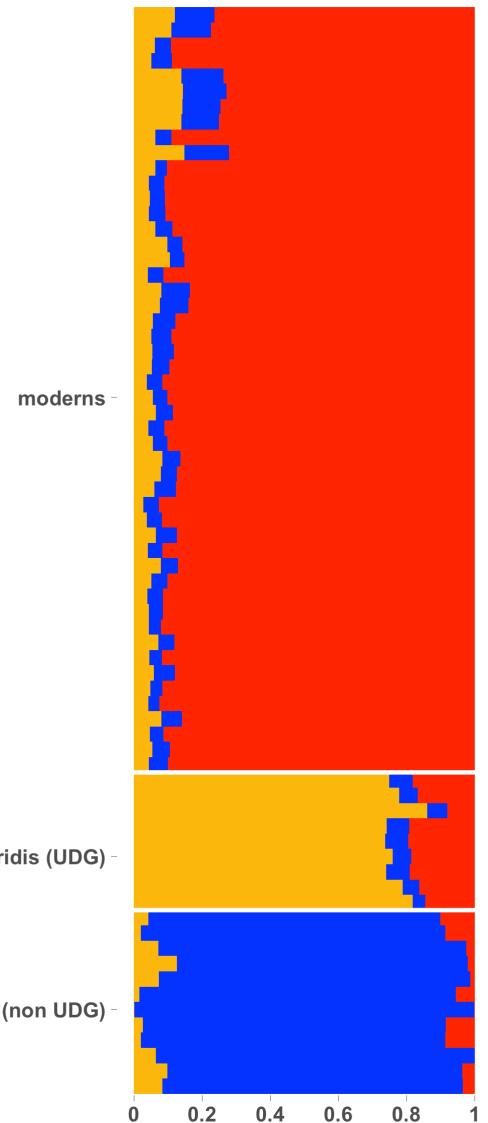
Applications of archaic

4 case studies

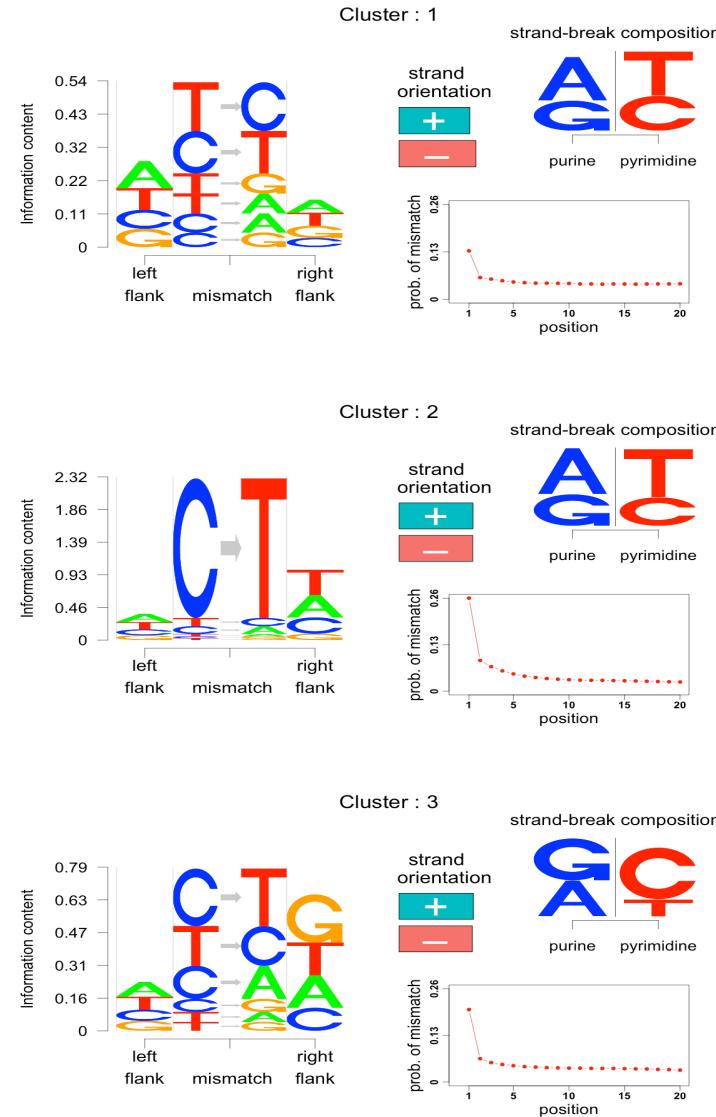
- Comparison of UDG treated aDNA (*Lazaridis+2014*), non UDG treated aDNA (Pinhasi et al) and moderns (1000G)
- Large scale aDNA due to *Mathieson+2015*
- *Lindo+2016* moderns and ancient DNA analysis
- Contamination study using simulations and *Mathieson+2015* data

Case Study 1 : 1000G moderns, Lazaridis and Pinhasi

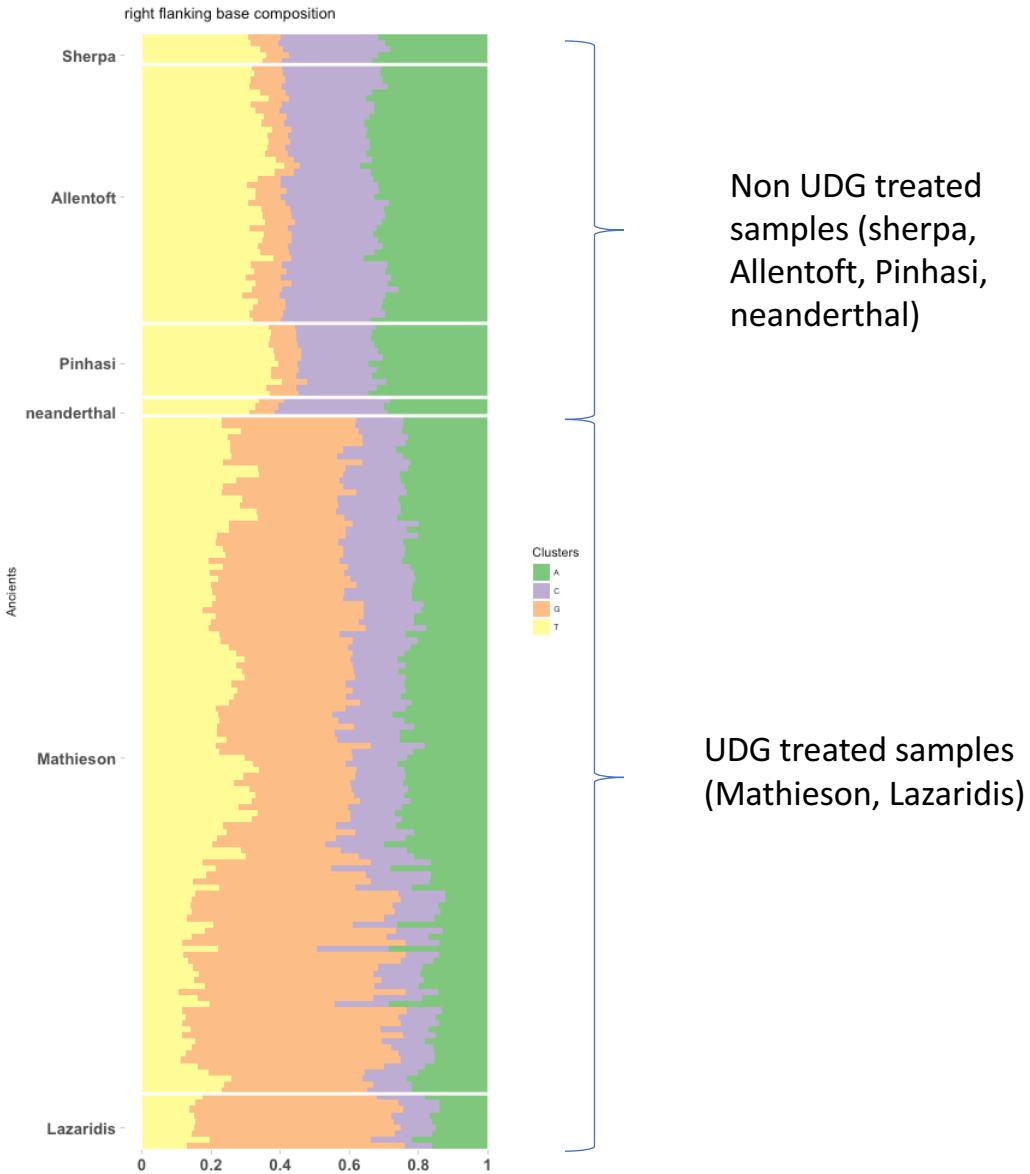
aRchaic pops

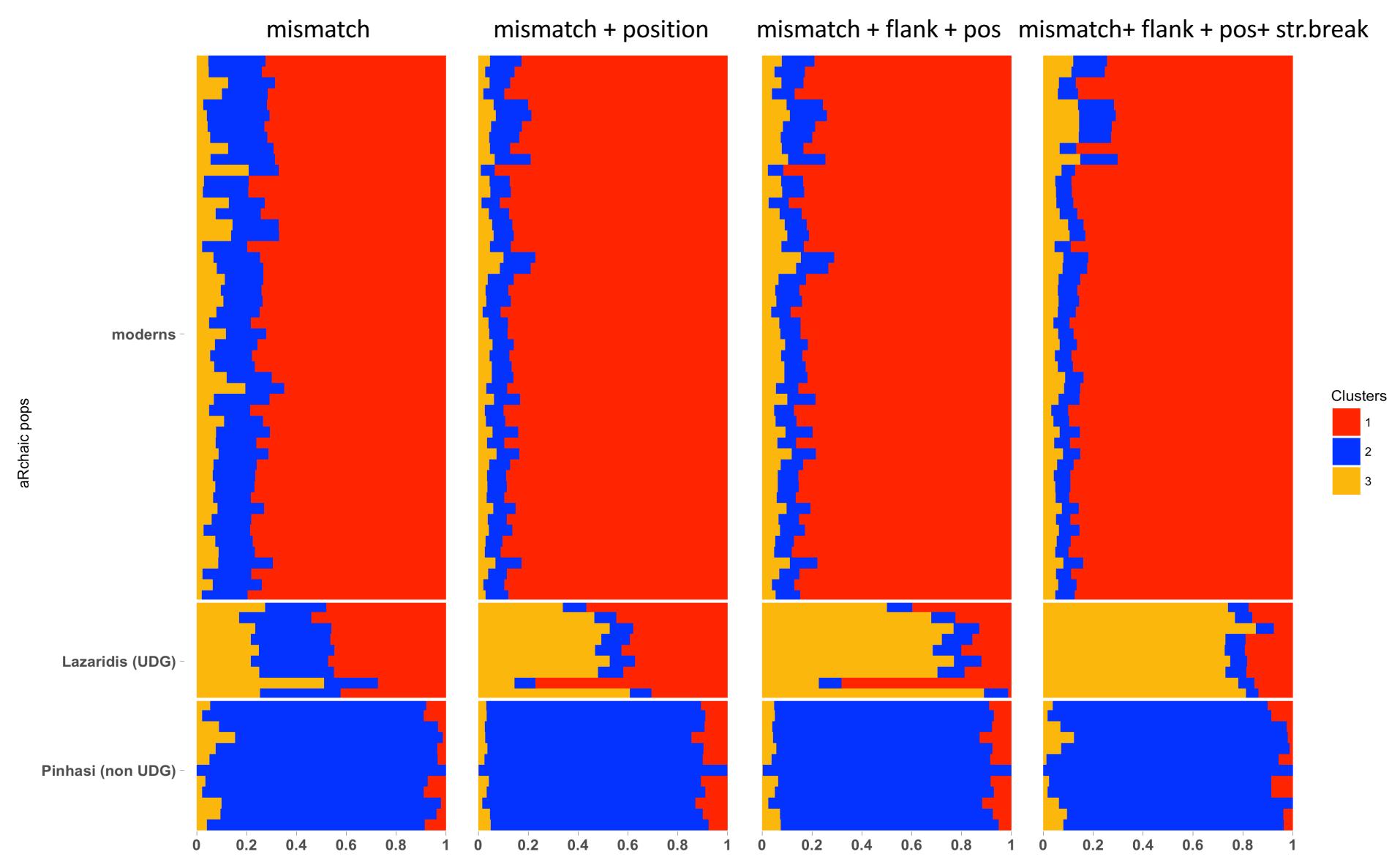


Clusters
1
2
3

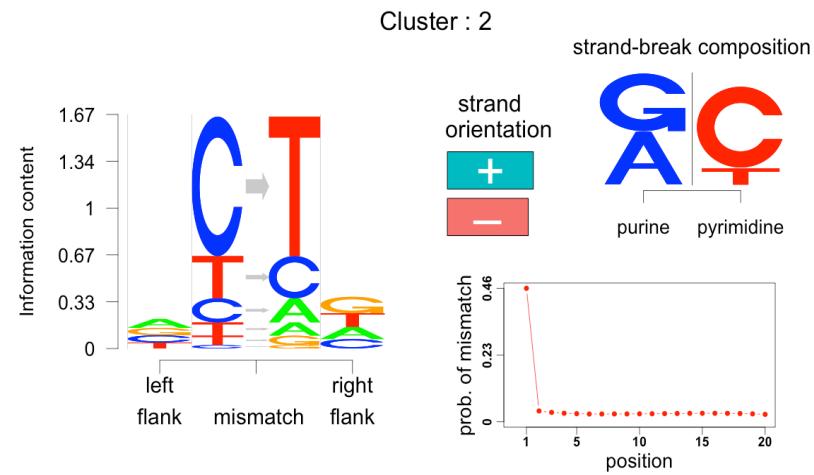
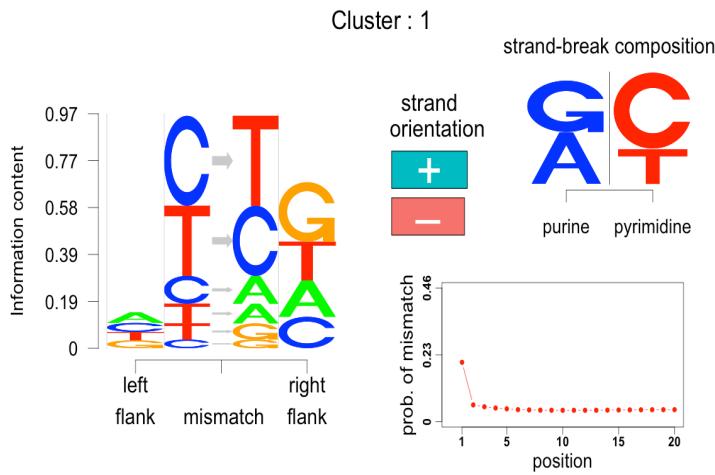
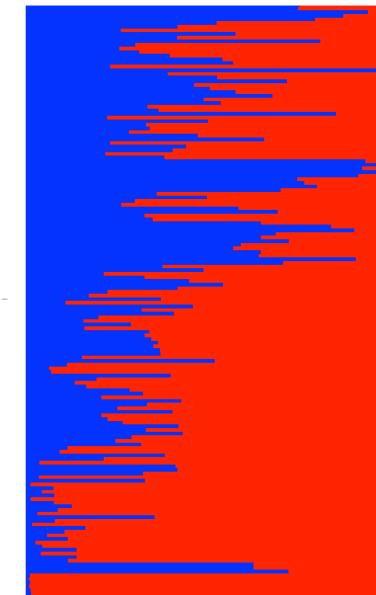
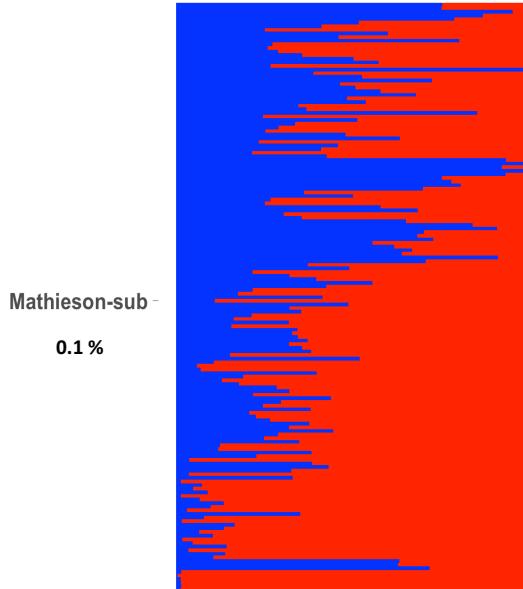


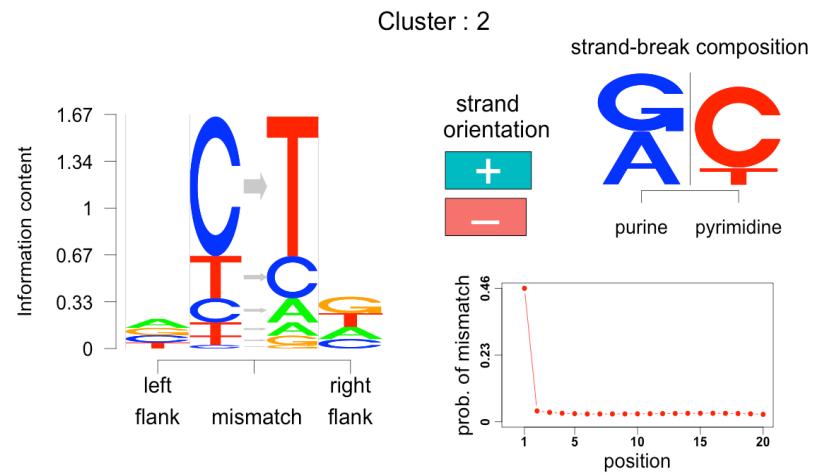
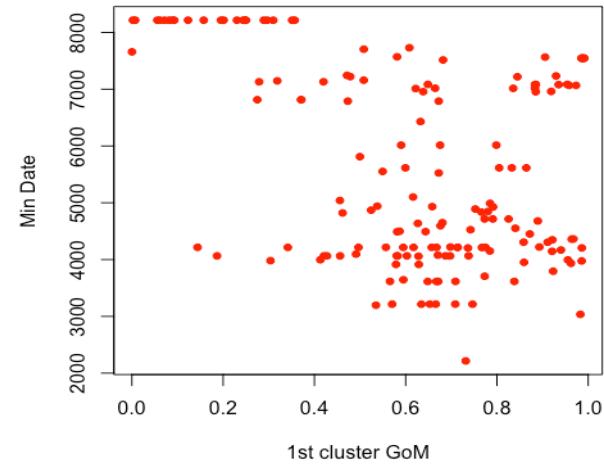
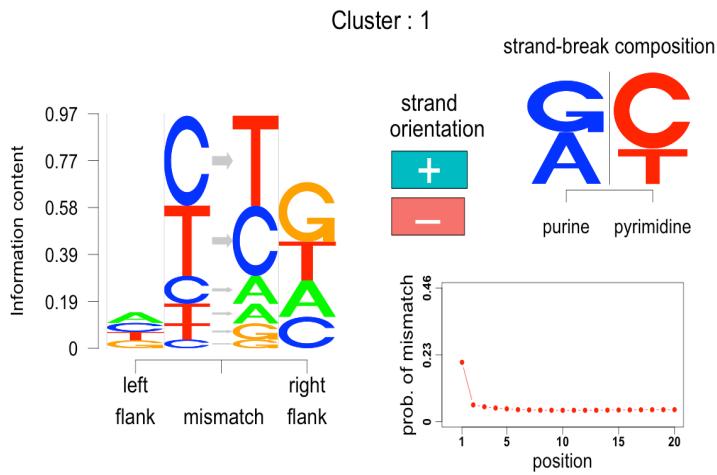
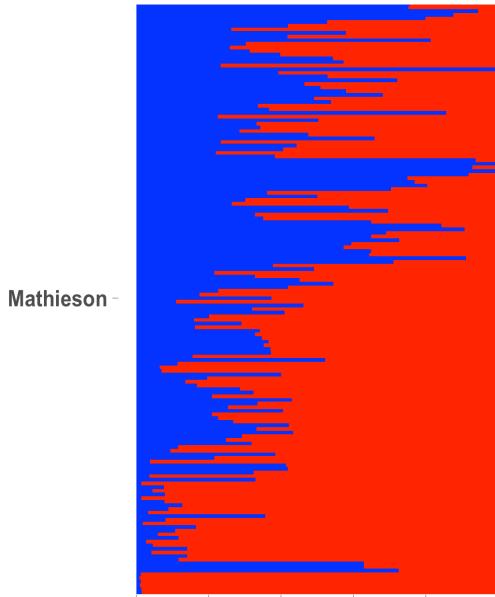
We look at the right flanking base composition for all mutations of type C->T occurring inside 10 bases from start of the read on the "+" strand.





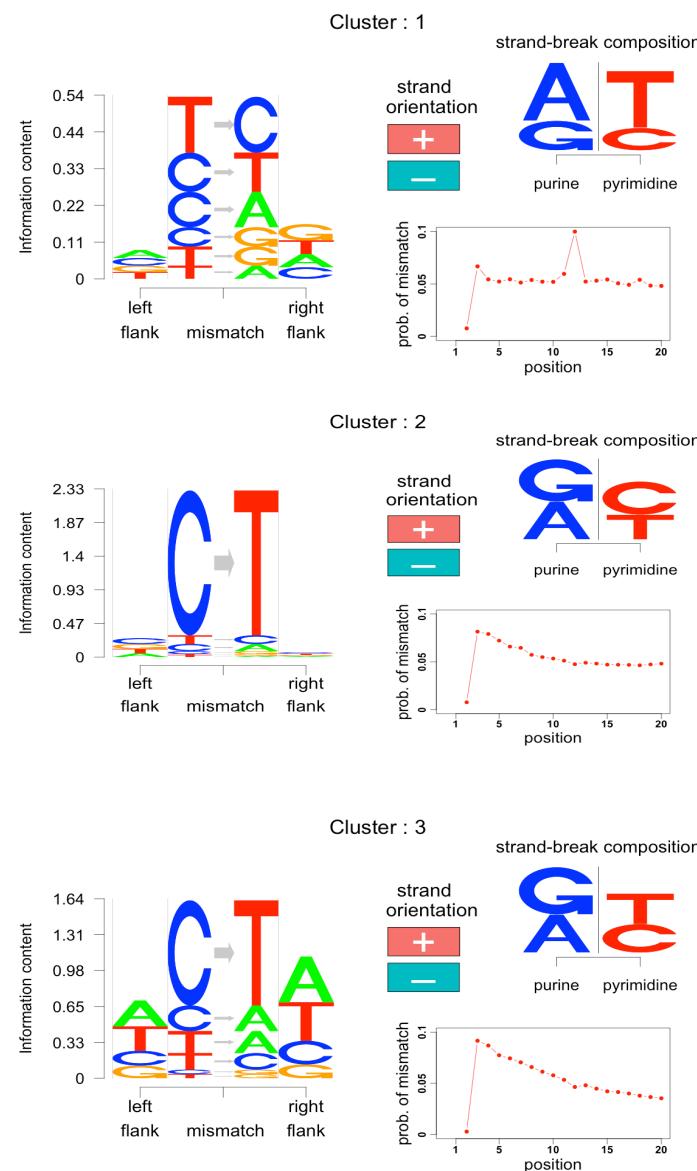
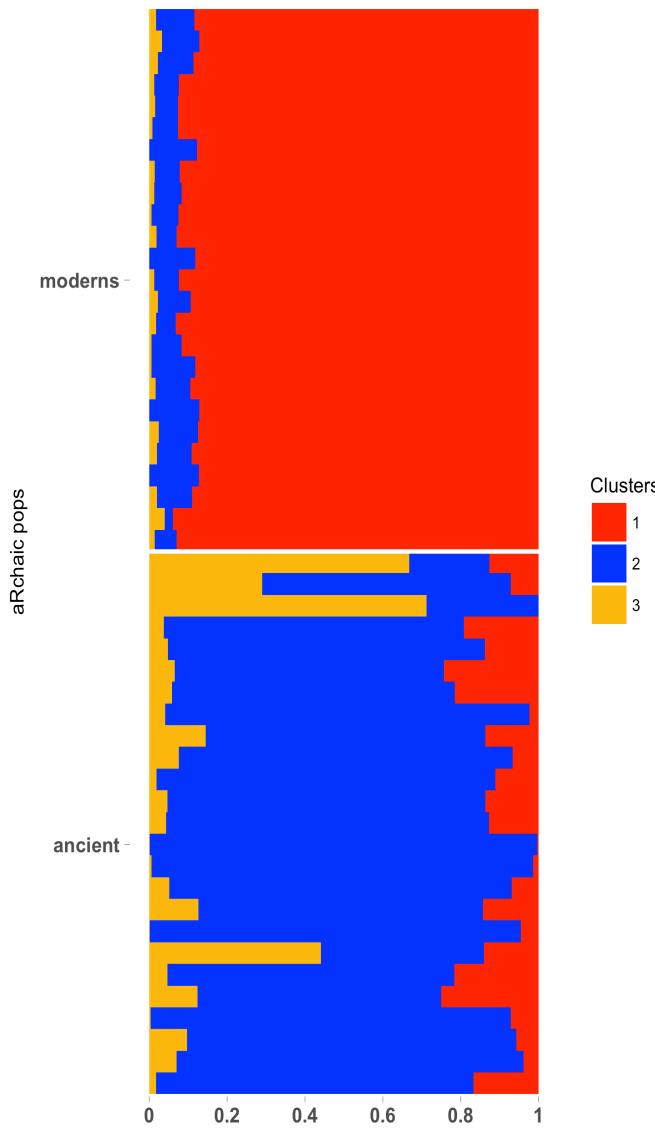
Case Study 2 : Mathieson data



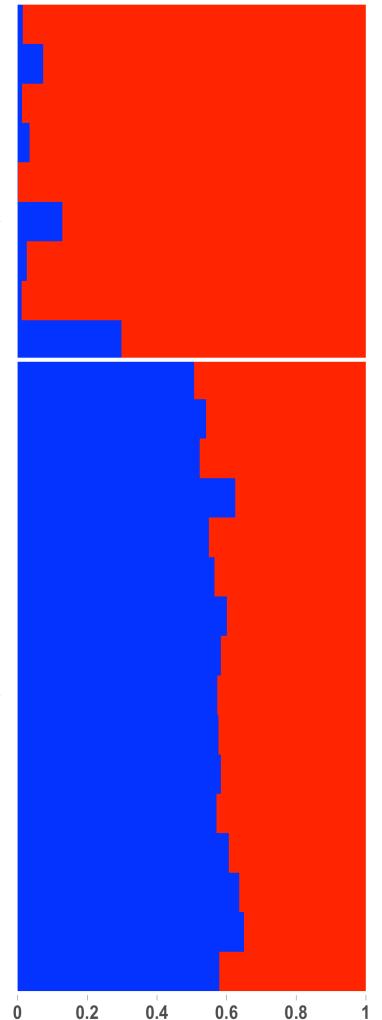


Case Study 3 : Lindo et al data

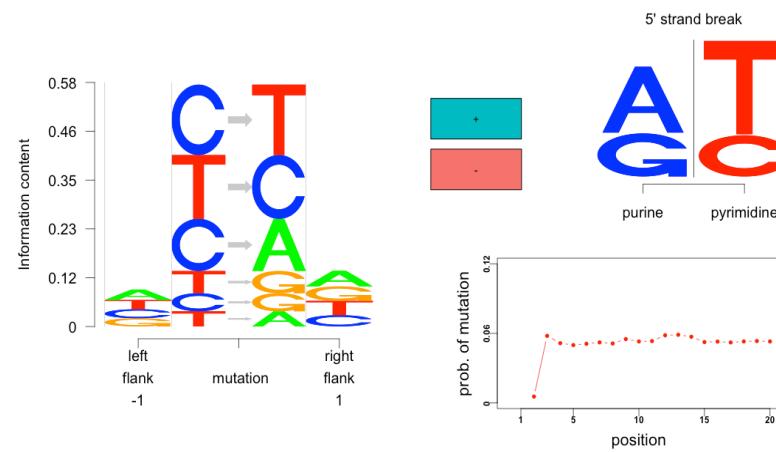
StructurePlot: K=3



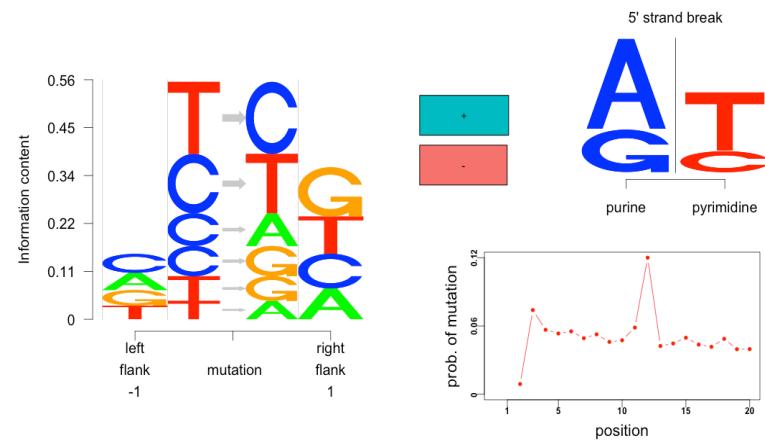
StructurePlot: K=2



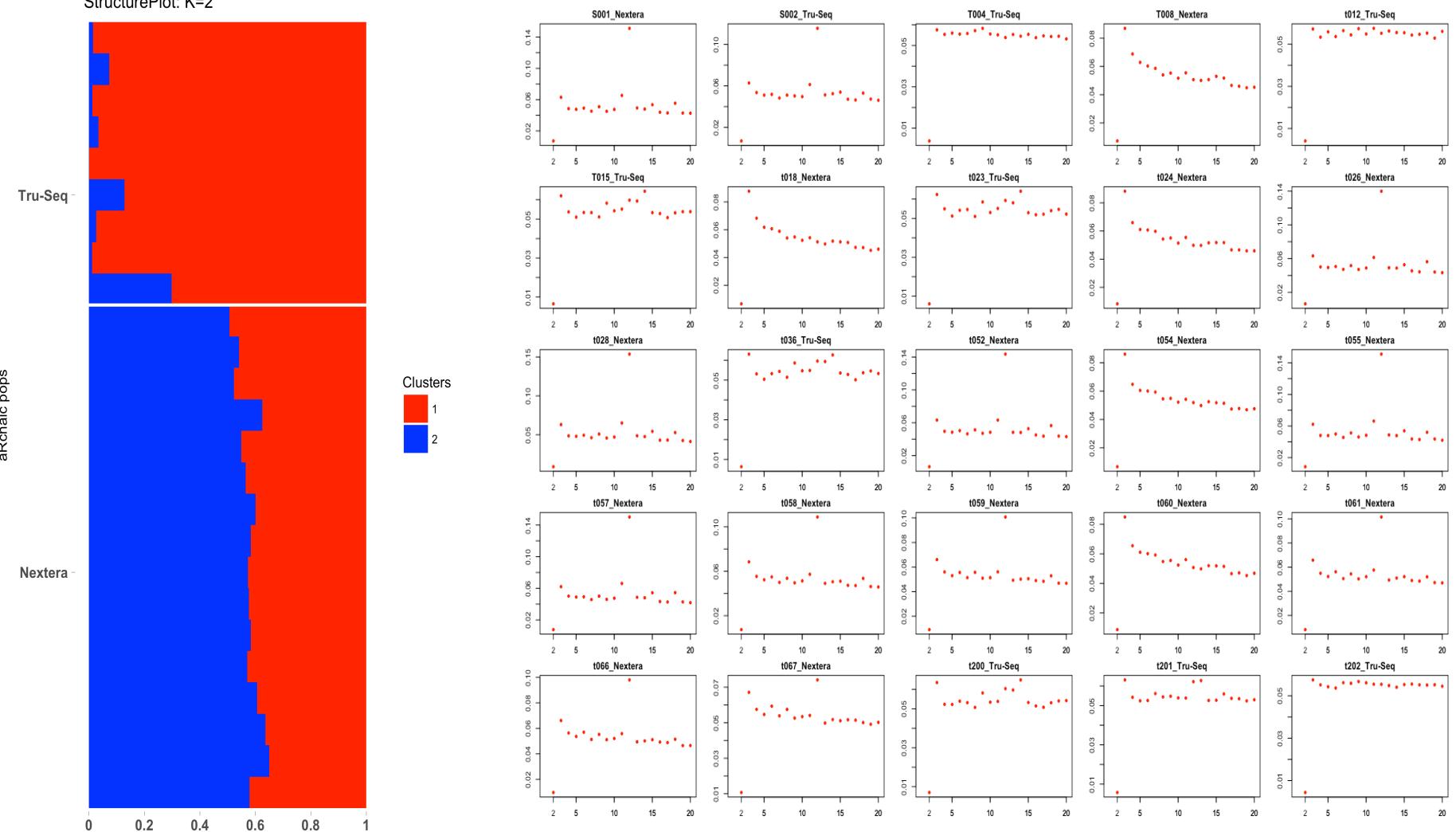
Cluster : 1



Cluster : 2



StructurePlot: K=2

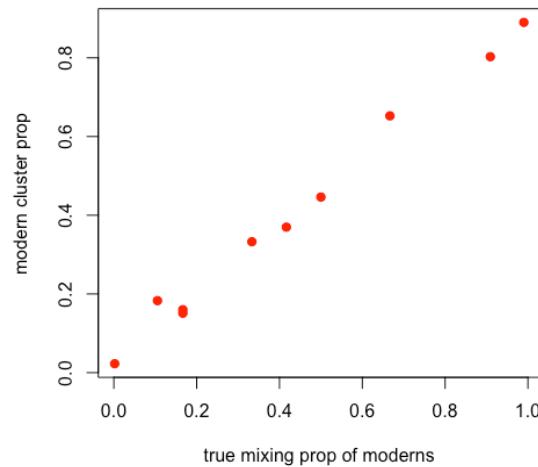


Case Study 4 : Contamination

Pool all mismatch patterns from a modern sample and all mismatch patterns from an ancient sample and mix them in a certain ratio (normalized also by total number of mutations).

Then run **aRchaic** on these mixed samples with other modern and ancient samples and see whether the grades of membership match with the mixing proportion.

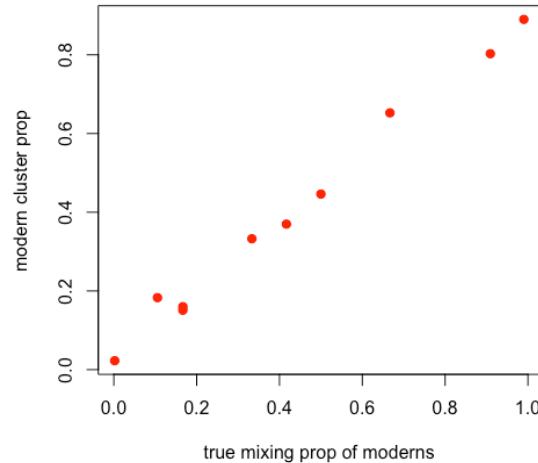
Pinhasi aDNA + moderns + their mixed samples



Pool all mismatch patterns from a modern sample and all mismatch patterns from an ancient sample and mix them in a certain ratio (normalized also by total number of mutations).

Then run **aRchaic** on these mixed samples with other modern and ancient samples and see whether the grades of membership match with the mixing proportion.

Pinhasi aDNA + moderns + their mixed samples

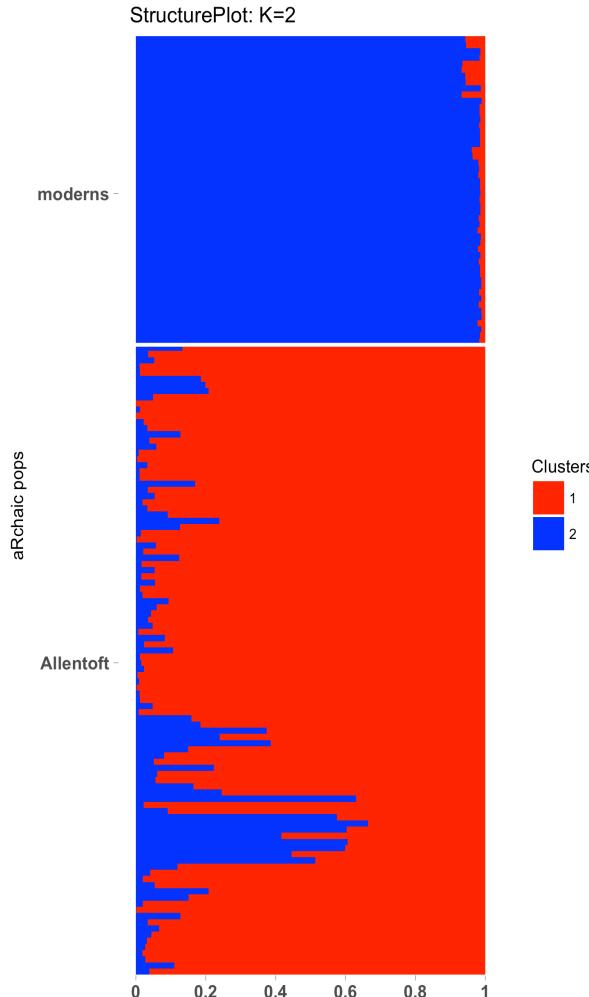


Issues with contamination:

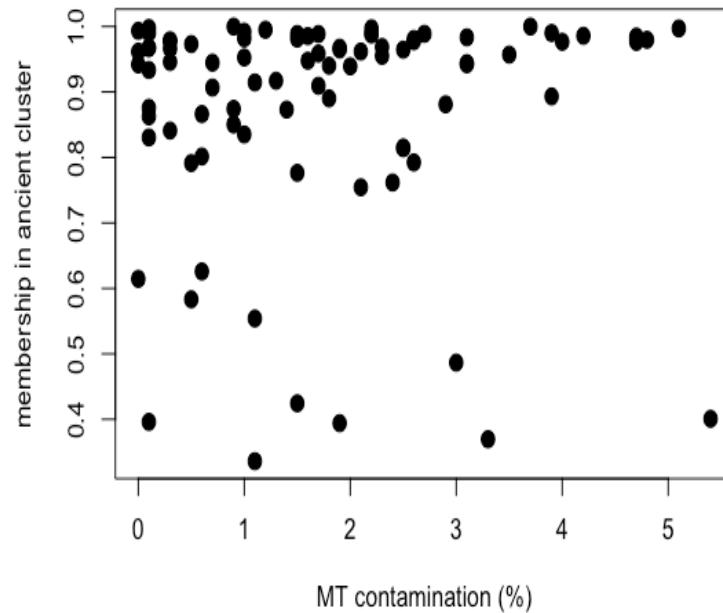
1) **aRchaic** only takes into account mismatches in calling mixing estimate and essentially throws away those reads that do not contain any mutations.
General contamination models take care of that information

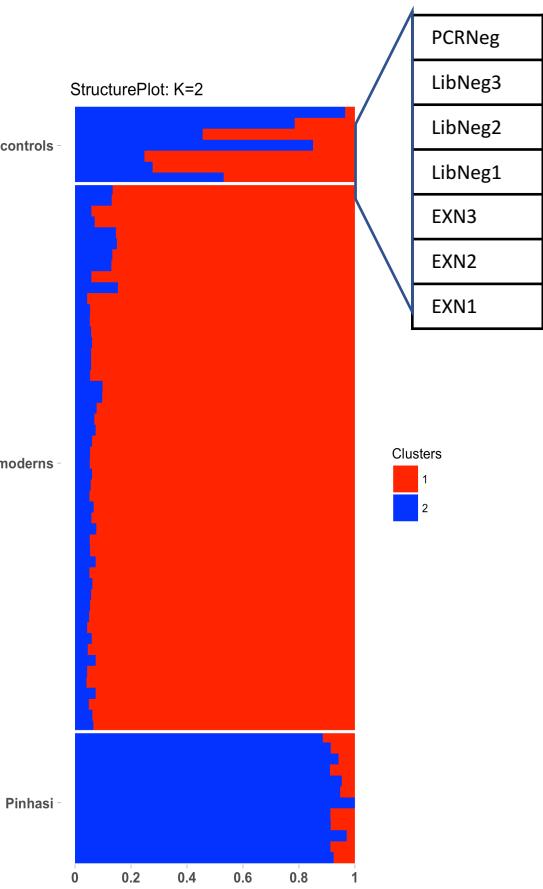
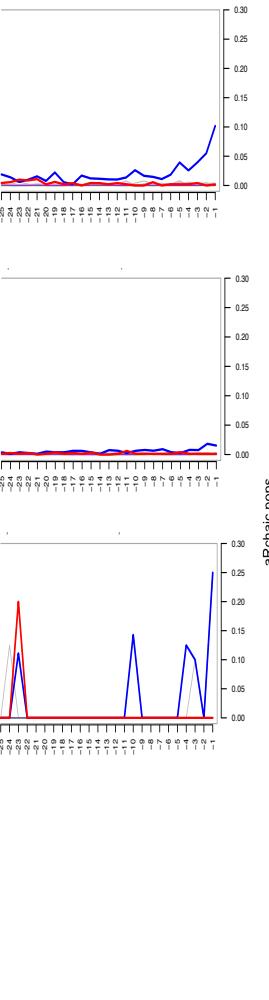
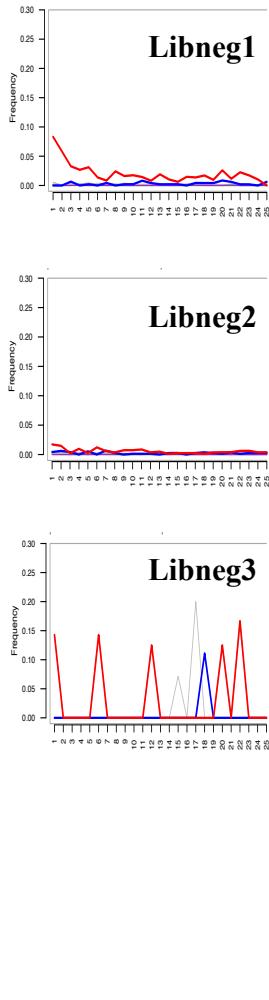
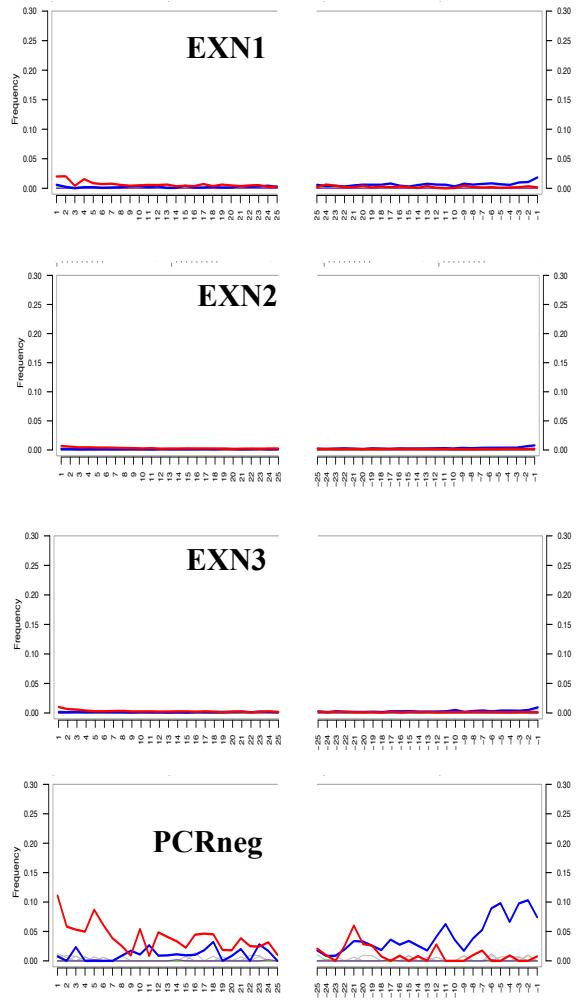
2) contemporary DNA that contaminates fossils also degrades to shorter length and accumulates mis-incorporations similar to ancient DNA (Briggs+2007, Malmström+2005, Sampietro+2006).

Moderns + Allentoft



mtDNA contamination estimate vs
aRchaic comparison for Allentoft
samples





Acknowledgements

Matthew Stephens
John Novembre
Anna Di Rienzo

David Witonsky
John Lindo
Anna Gosling
Choongwong Jeong

John Blischak
Peter Carbonetto
Joe Marcus
Yuichi Shiraishi

Software

aRchaic is available as R package through
devtools::install_github("kkdey/aRchaic")

Webpage (still under work):
kkdey.github.io/aRchaic

