

# Simulation of ancient DNA sequences using transformer-based techniques.

Théo Boury<sup>1,2</sup>, Jazeps Medina-Tretmanis<sup>3</sup>, Maria Avila-Arcos<sup>4</sup>,  
Emilia Huerta-Sanchez<sup>3</sup>, Burak Yelmen<sup>2,5</sup>, Flora Jay<sup>2</sup>

<sup>1</sup>Computer Science Department, Ecole Normale Supérieure de Lyon, France; <sup>2</sup>U Paris-Saclay, CNRS, INRIA, Laboratoire Interdisciplinaire des Sciences du Numérique, France; <sup>3</sup> Center for Computational Molecular Biology, Brown U, USA; <sup>4</sup> International Laboratory for Human Genome Research, U Nacional Autónoma de México, México; <sup>5</sup>Institute of Genomics, U of Tartu, Estonia.

# Ancient DNA specificities

Undamaged DNA

A T C G T

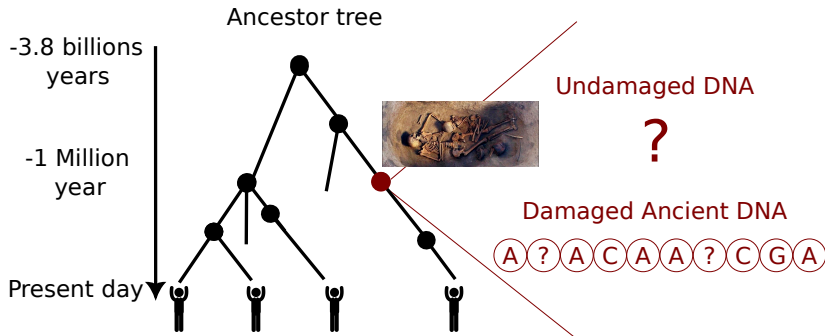
Damaged Ancient DNA

A ? A C A A ? C G A

Difficulties with Ancient DNA (aDNA):

- | Degrades over time
- | Contaminated by external DNA
- | More missing data and errors than modern DNA

# Why study Damaged Ancient DNA?



## ARTICLE

### Ancient gene flow from early modern humans into Eastern Neanderthals

Marta Sabido<sup>1</sup>, Dan Groux<sup>1</sup>, Mathieu Hublin<sup>1</sup>, Christophe Filippucci<sup>1</sup>, Jean-François Meunier<sup>1</sup>, Martin Elchert<sup>1</sup>, Quentin Pons<sup>1</sup>, Françoise Audouze<sup>1</sup>, Céline Lahr<sup>1</sup>, Françoise Bon<sup>1</sup>, Marco de la Haye<sup>1</sup>, Antoine Lécuyer<sup>1</sup>, Pierre-André Crochet<sup>1</sup>, Fabrice Lebrun<sup>1</sup>, Jean-Louis Lhomme<sup>1</sup>, Thomas G. Shved<sup>1</sup>, Julia M. Aulic<sup>1</sup>, Olivier Heuvelink<sup>1</sup>, Steffen Pääbo<sup>1</sup>, Martin Wehner<sup>1</sup>, Adam Izuel<sup>1</sup> & Sébastien Caillaud<sup>1</sup>

SCIENCE ADVANCES | RESEARCH ARTICLE

EVOLUTIONARY BIOLOGY

#### Genetic ancestry changes in Stone to Bronze Age transition in the East European plain

Lehti Saag<sup>1</sup>, Sergey V. Vasylyev<sup>2</sup>, Liwei Yan<sup>1</sup>, Natalia V. Kosonikova<sup>3</sup>, Dmitri V. Gerasimov<sup>4</sup>, Svetlana V. Ostikhina<sup>5</sup>, Samuel J. Griffith<sup>6</sup>, Aina Sotnik<sup>7</sup>, Levent Saag<sup>1</sup>, Eugenia D'Elia<sup>8</sup>, Ene Metspalu<sup>1</sup>, Maarek Reidla<sup>1</sup>, Siiri Nootsi<sup>1</sup>, Toomas Kivisild<sup>1,9</sup>, Christiana Lyn Schibler<sup>1,9</sup>, Kristina Tambets<sup>1</sup>, Alvar Kirisik<sup>10</sup>, Mait Metspalu<sup>1</sup>

## Article

### Large-scale migration into Britain during the Middle to Late Bronze Age

<https://doi.org/10.1093/nar/nkz104>

Received: 20 December 2019

Accepted: 20 November 2020

Published online: 22 December 2020

All of our articles and their associated metadata are freely available for re-use under a CC-BY license.

For more information on this article, please visit the article page on the NAR website.

Early European Farmers (EEF) chromosome of the Early Bronze Age. Technology

like, we generated genome-wide data from 70 individuals, including high-coverage

Middle to Late Bronze and Iron Age individuals by 12-fold and 100-fold coverage

Cell

#### The genomic origins of the world's first farmers

## RESEARCH ARTICLE SUMMARY

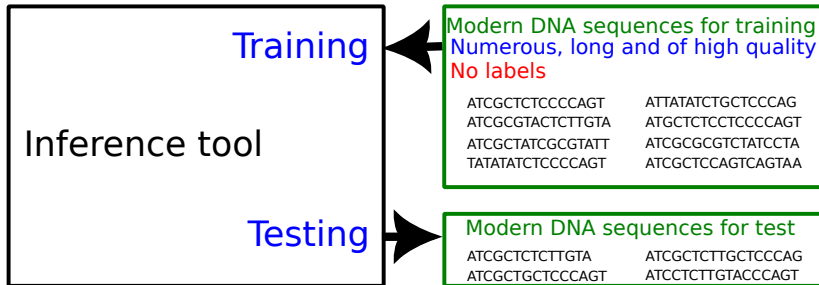
HUMAN EVOLUTION

### The formation of human populations in South and Central Asia

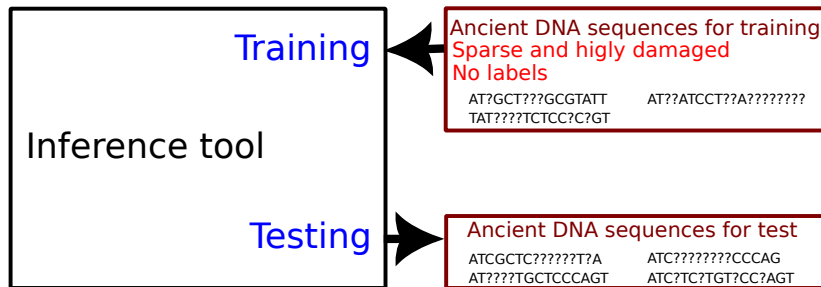
Preprint 20. November 2019

Article

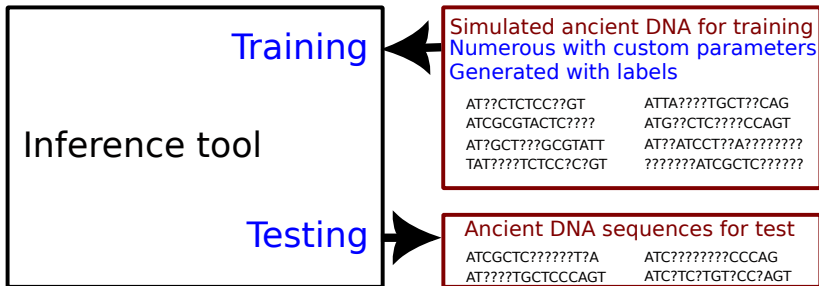
# DNA sequences for inference purposes....



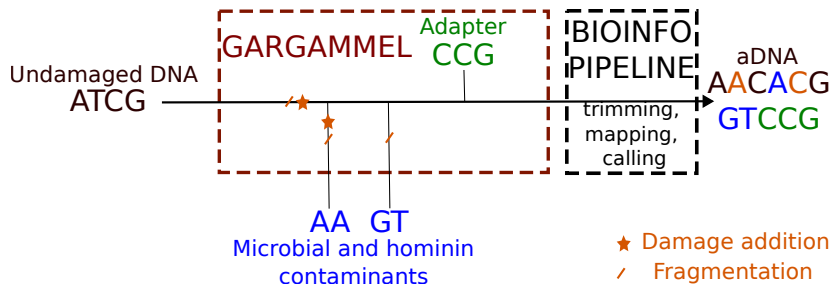
# aDNA sequences for inference purposes....



... required simulation of aDNA sequences



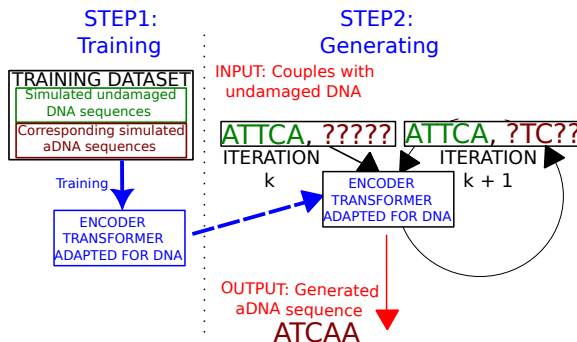
# State-of-the-art aDNA simulator: Gargammel<sup>1</sup>



- | Gargammel complexity:  $O(n \times c \times f)$
- | With  $c$ , the desired coverage and  $f$ , the number of fragments sampled by Gargammel
- |  $f$  can lead to a large overhead in practice

<sup>1</sup>Renaud et al, 2016, Bioinformatics

# Achieved result: our new seq-to-seq aDNA simulator



- | **Generate** aDNA sequences from undamaged ones
- | **Method**: Iterative process over a specific encoder-only transformer
- | **Data simulation**: Undamaged sequences: Msprime<sup>2</sup>.  
Damaged sequences: pipeline around Gargammel<sup>3</sup>
- | **Generation complexity**:  $O(n^3)$ , of interests compared to Gargammel

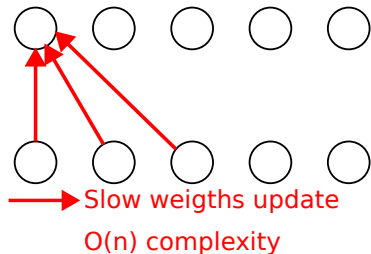
<sup>2</sup>Baumdicker et al, 2021, Genetics

<sup>3</sup>Jazeps et al, 2023, ?

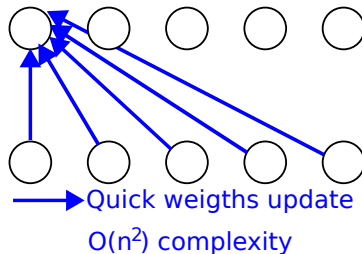


# Attention interest<sup>4</sup> versus convolution

## Convolution



## Attention



---

<sup>4</sup>Vaswani et al, 2017

# Our use of the pretrained DNABERT model<sup>5</sup>

I Max seq. length: 512 nucleotides. Comp.: 12 Transformer layers

<sup>5</sup>Ji et al, 2021, Bioinformatics

# Our data and our re-tunings

- | A **mask prediction task** to predict the aDNA part
- | A **binary classification task** to measure if the aDNA part is a plausible "translation" of the undamaged DNA

# Generation algorithm using Mask Prediction <sup>6</sup>only

- | Complexity with Mask Prediction alone  $\mathcal{O}(n^3)$
- | With  $n$ , size of the input

<sup>6</sup>Inspired by bert-gen, Wang et al, 2019

# Use of classification as a complement for mask prediction in the generation

- | Classification: Is the aDNA sequence a "plausible" translation for the undamaged sequence?
- | Complexity with the K-Top table:  $O(Kn^3)$

# Results

- | We do alignments<sup>7</sup> with **30 chunks of aDNA**
- | Aligned sequences are identical at **74 percent in average**
- | We count similar positions to the exclusion of gaps and missing data

<sup>7</sup>Wheeler et al, 2000, Nucleic Acids Research

# Perspective: push further complexity and performances

- | First leverage: **Sparse<sup>8</sup> or linear<sup>9</sup> attention** instead of full attention  
Reduce attention to  $O(n)$  instead of  $O(n^2)$ .
- | Second leverage: **Use of SNPs**  
Counteracts the 512 nucleotide limitation and "diminishes"

---

<sup>8</sup>Zaheer et al, 2021

<sup>9</sup>Nesterenko et al, 2022

# Conclusion and future work

- | A new **seq-to-seq simulation technique** for ancient DNA sequences
- | **Complexity in  $O(n^3)$**  in simpler case, complemented to the use of batches in practice

## Future work:

- | Define new criteria to **assess the quality of sequences**
- | Do our own pre-training instead of using DNABERT's one.
- | Define an **"in-between" classification task** when using k-Top table
- | General advance in transformers will be at our advantage:  
Integration of sparse attention



# Acknowledgements

## Contributors:

- | Jazeps Medina-Tretmanis
- | Maria Avila-Arcos
- | Emilia Huerta-Sanchez
- | Antoine Szatkownik (\Generative model for genomics" poster!)
- | Burak Yelmen
- | Flora Jay
- | Guillaume Charpiat

## Laboratories and funding agencies:

# Our ancient DNA simulation using encoder-transformer

