

Simulation of ancient DNA sequences using transformer-based techniques.

Théo Boury^{1,2}, Jazeps Medina-Tretmanis³, Maria Avila-Arcos⁴,
Emilia Huerta-Sanchez³, Burak Yelmen^{2,5}, Flora Jay²

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

Ancient DNA specificities

Undamaged DNA



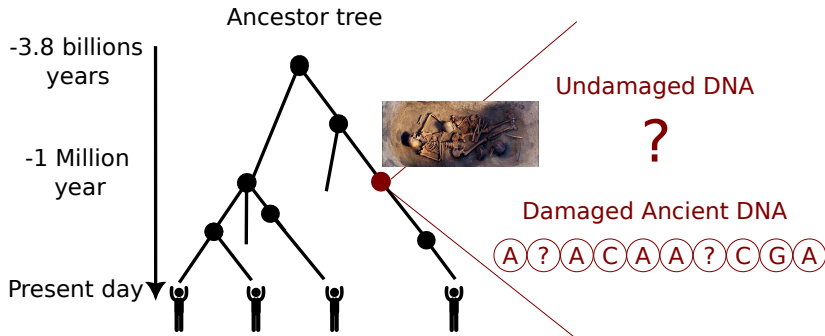
Damaged Ancient DNA



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 Sabbe¹, Dan Gascón¹, Mélanie Hublin¹, Christophe Filippucci¹, Javier Prado-Martinez¹, Martin Elstner¹, Quamrul Hossain¹, Jingxiu A. Bao¹, Carlos Lalueza-Fox¹, Mariona Estroix de Serra¹, Prasad Reddy¹, Steffen Richter¹, Tobias Meyer¹, Irene Galardi¹, Teresa Gilibert¹, Rasmus Nielsen¹, John M. Aitken¹, Svante Pääbo¹, Martin Nieber¹, Adam Siegel¹ & Ségolène Coudane

SCIENCE ADVANCES | RESEARCH ARTICLE

EVOLUTIONARY BIOLOGY

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

Lehti Saag¹, Sergey V. Vasylyev², Liwei Yan¹, Natalia V. Kosonikova³, Dmitri V. Gerasimov⁴, Svetlana V. Ostikhina⁵, Samuel J. Griffith⁶, Aina Sotnik⁷, Levent Sag⁸, Eugenia D'Elia⁹, Ene Metspalu¹, Maere Heida¹, Sifri Nootsi¹, Toomas Kivisild¹⁰, Christiana Lyn Schibler¹⁰, Kristina Tambets¹, Alvar Kirisik¹¹, Mait Metspalu¹

Article

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

<https://doi.org/10.1093/aag/abaa014>

Received: 20 December 2020

Accepted: 20 November 2021

Published online: 22 December 2021

All of our articles and their associated supplements are licensed under a Creative Commons Attribution 4.0 International License.

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

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

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

Cell

The genomic origins of the world's first farmers

RESEARCH ARTICLE SUMMARY

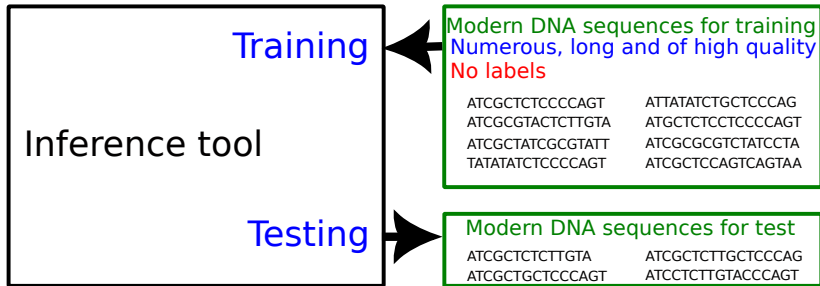
HUMAN EVOLUTION

The formation of human populations in South and Central Asia

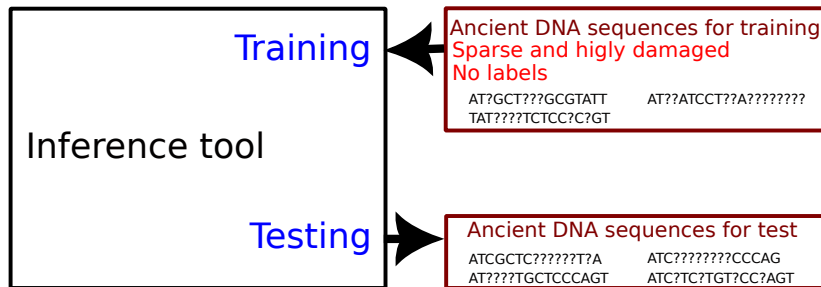
Preprint 20. November 2021

Article

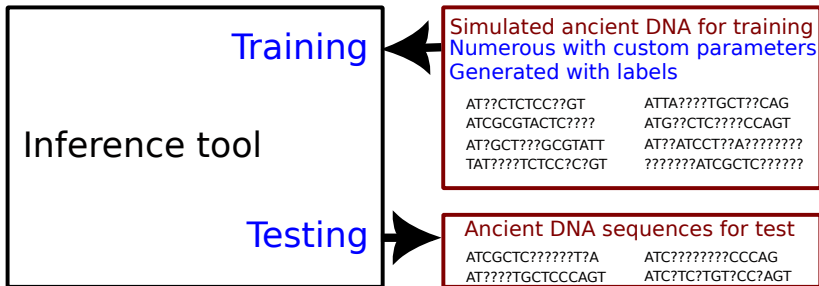
DNA sequences for inference purposes....



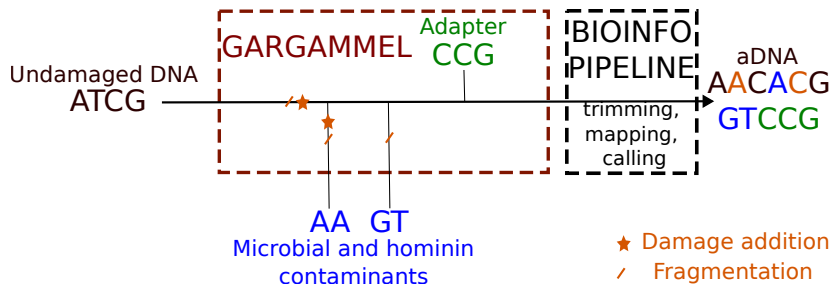
aDNA sequences for inference purposes....



... required simulation of aDNA sequences



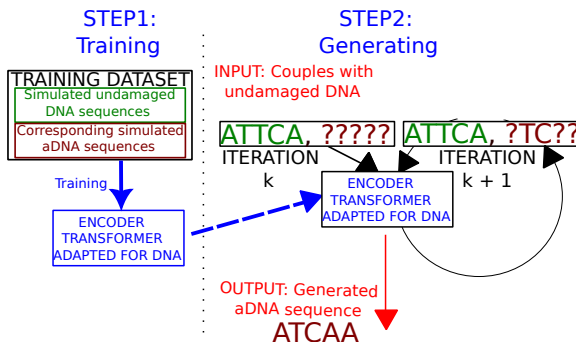
State-of-the-art aDNA simulator: Gargammel¹



- ▶ 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

¹Renaud et al, 2016, Bioinformatics

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



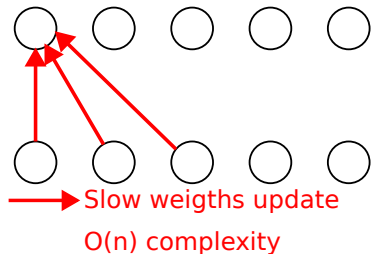
- ▶ **Generate** aDNA sequences from undamaged ones
- ▶ **Method**: Iterative process over a specific encoder-only transformer
- ▶ **Data simulation**: Undamaged sequences: Msprime². Damaged sequences: pipeline around Gargammel³
- ▶ **Generation complexity**: $O(n^3)$, of interests compared to Gargammel

²Baumdicker et al, 2021, Genetics

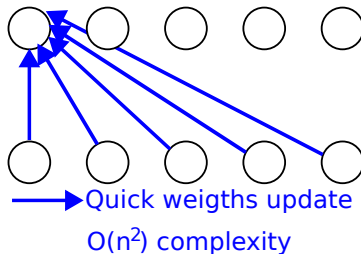
³Jazeps et al, 2023, ?

Attention interest⁴ versus convolution

Convolution

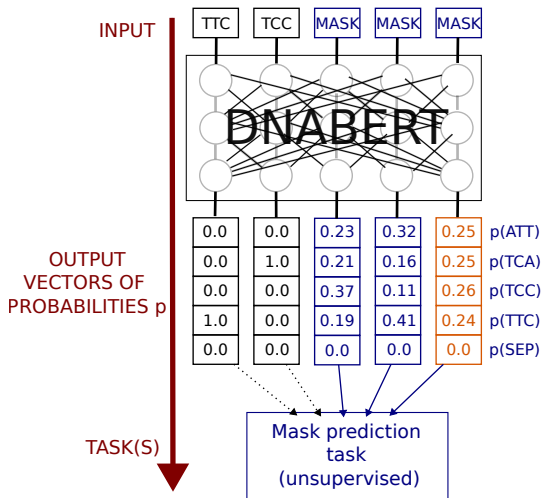


Attention



⁴Vaswani et al, 2017

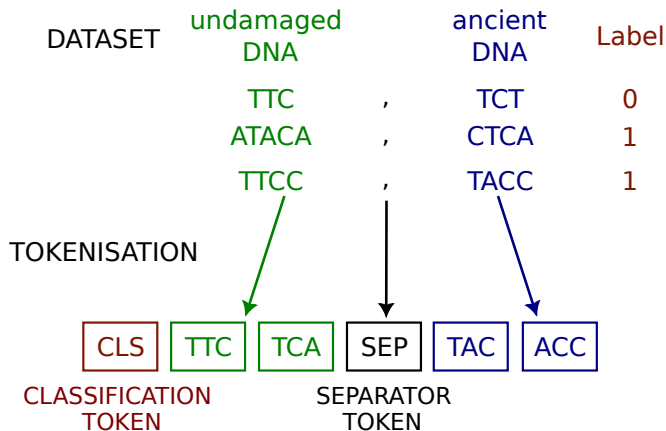
Our use of the pretrained DNABERT model⁵



- ▶ Max seq. length: 512 nucleotides. Comp.: 12 Transformer layers

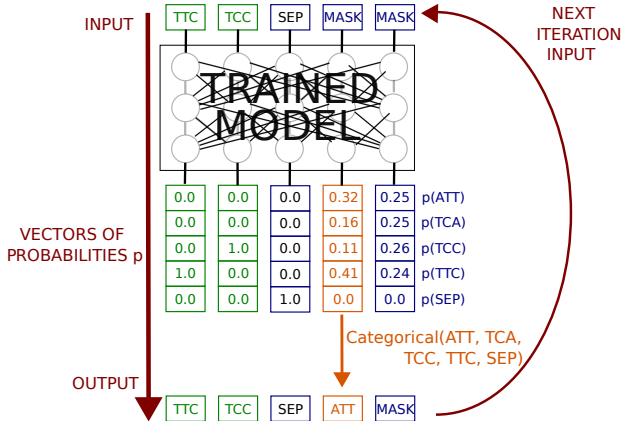
⁵Ji et al, 2021, Bioinformatics

Our data and our fine-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 only⁶



- ▶ Complexity with Mask Prediction alone: $O(n^3)$
- ▶ With n , size of the input

⁶Inspired by bert-gen, Wang et al, 2019

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

$P(\text{CLS}=1)$	3-TOP aDNA sequences at iteration k				$P(\text{CLS}=1)$ at iteration $k + 2$	Candidate aDNA sequence at iteration $k + 1$			
0.32	MASK	MASK	TCA	MASK	0.34	MASK	ATT	ATT	MASK
0.37	MASK	ATT	MASK	MASK					
0.33	TTC	TCC	MASK	MASK					

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

Perspective: push further complexity and performances

- ▶ **First leverage:** Sparse⁸ or linear⁹ 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" n

	Full sequences	Intermediates	SNPs
	A T T A G G A C G	A A G G	0 1 1 0
	A T T A C G A C A	A A C A	0 1 0 1
	C T T C G G A C G	C C G G	1 0 1 0
Positions	1 2 3 4 5 6 7 8 9	1 4 5 9	1 4 5 9

⁸Zaheer et al, 2021

⁹Nesterenko 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 a 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

YOU!

Laboratories and funding agencies:



Our ancient DNA simulation using encoder-transformer

