



HAL
open science

Topological data analysis and multiple kernel learning for species identification of modern and archaeological small ruminants

Manon Vuillien, Davide Adamo, Emmanuelle Vila, Agraw Amane, Thierry Argant, Daniel Helmer, Marjan Mashkour, Abdelkader Mousous, Olivier Notter, Elena Rossoni-Notter, et al.

► **To cite this version:**

Manon Vuillien, Davide Adamo, Emmanuelle Vila, Agraw Amane, Thierry Argant, et al.. Topological data analysis and multiple kernel learning for species identification of modern and archaeological small ruminants. 2024. hal-04779367v2

HAL Id: hal-04779367

<https://hal.science/hal-04779367v2>

Preprint submitted on 7 Feb 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

Topological data analysis and multiple kernel learning for species identification of modern and archaeological small ruminants

Manon Vuillien^a, Davide Adamo^{a,b}, Emmanuelle Vila^c, Agraw Amane^{d,e}, Thierry Argant^f, Daniel Helmer^c, Marjan Mashkour^g, Abdelkader Moussous^h, Olivier Notter^h, Elena Rossoni-Notter^h, Isabelle Théry^a, Marco Corneli^{a,b}

^a Université Côte d'Azur, UMR 7264 CEPAM, CNRS, Nice, France

^b Université Côte d'Azur, Inria, CNRS, Laboratoire J.A. Dieudonné, Maasai team, Nice, France

^c Université Lumière Lyon II, UMR 5133 Archéorient, CNRS, Lyon, France

^d Department of Microbial Cellular and Molecular Biology, Addis Ababa University, Ethiopia

^e ILRI Livestock Genetics Program, Addis Ababa, Ethiopia

^f Eveha études et valorisation archéologique, UMR 5138 ArAr, Lyon, France

^g AASPE Muséum national d'Histoire naturelle, CNRS, Paris, France

^h Musée d'Anthropologie préhistorique de Monaco, Monaco

Abstract

The faunal remains from numerous Holocene archaeological sites across southwest Asia frequently include the bones of various wild and domestic ungulates, such as sheep, goats, ibexes, roe deer and gazelles. These assemblages may provide insight into hunting and animal husbandry strategies and offer palaeoecological information on ancient human societies. However, the skeletons of these taxa are highly similar in appearance, which presents a challenge for accurate identification based on their bones. This paper presents a case study to test the potential of topological data analysis (TDA) and multiple kernel learning (MKL) for inter-specific identification of 150 3D astragali belonging to modern and archaeological specimens. The joint application of TDA and MKL demonstrated remarkable efficacy in accurately identifying wild species, with a correct identification rate of approximately 90%. In contrast, the identification of domestic species exhibited a lower success rate, at approximately 60%. This low rate of identification of sheep and goat species is attributed to the morphological variability of domestic breeds. Moreover, while these methods assist in clearly identifying wild taxa from one another, they also highlight their morphological diversity. In this context, TDA and MKL could be invaluable for investigating intra-specific variability in domestic and wild animals. These methods offer a means of expanding our understanding of past domestic animal selection practices and techniques. They also facilitate an investigation into the morphological evolution of wild animal populations over time.

Keywords

Machine learning, Osteology, Herbivores, Zooarchaeology, Topological data analysis, Multiple kernel learning

1. Introduction

Small wild and domestic ungulates are frequently found in Holocene archaeological faunas and are likely to be found together in some contexts. For example, in the Near East, wild and domestic *Caprinae*

1 (domestic or wild goat and sheep, ibex), roe deer, and gazelle have geographical distributions whose
2 limits in ancient times are poorly defined and partly overlap for specific taxa (Uerpmann 1987). Each
3 of these species shows adaptation to a particular ecological habitat and specific plant resources. Roe
4 deer prefers areas of mixed forest and grassland. Ibex is a mixed feeder (browser and grazer) living in
5 mountainous regions. Gazelle occurs in waterless steppe, semi-arid, and desert environments. These
6 ungulates provide information on hunting and husbandry strategies and palaeoecological information
7 on the climate and environment of ancient societies, the diversity of natural habitats, and regional
8 variations in terms of aridification or vegetation cover capacity, agriculture, and deforestation (Tsahar
9 *et al.* 2009). However, the skeletons of these taxa are very close morphologically, which poses a
10 problem for specific identification based on their bones, as evidenced by numerous methodological
11 studies conducted on these species (Fernandez 2001; Salvagno & Albarella 2017; Sipilä *et al.* 2023;
12 Zeder & Lapham 2010; Zeder & Pilaar 2010) over the past fifty years. Furthermore, specific taxa, such
13 as gazelles, face challenges in distinguishing between species due to the absence of anatomical criteria
14 (Buitenhuis 1988; Gudea & Stan 2012; Peters 1989) or identifying sexual dimorphism (Munro *et al.*
15 2011) due to their morphological similarity.

16
17 Indeed, the taxonomic identification of remains of morphologically related species found in
18 archaeological contexts represents one of the key challenges that zooarchaeologists face. Traditionally,
19 the process of identifying bones or dental remains in archaeology is based on anatomical,
20 morphological and biometric criteria. These taxa are compared with their modern or fossil
21 counterparts documented in modern comparative osteological collections or represented in
22 anatomical atlases. It is often the case that identification criteria are provided by the literature which
23 have been tested and validated on large reference collections. Therefore, these criteria are very likely
24 to be accurate. However, this process can still be challenging due to factors such morphological
25 convergence within closely related species, potential absence of diagnostic criteria, intermediate
26 morphological characteristics, and intra-individual variability. Consequently, the use of anatomical
27 criteria available in the literature to differentiate these species is not always sufficient.

28
29 In recent decades, palaeogenetic (Alberto *et al.* 2018; Daly *et al.* 2018; Larsson *et al.* 2024; Lv *et al.*
30 2022) and palaeoproteomic (Fabrizi *et al.* 2024; Le Meillour *et al.* 2023; Pilaar Birch *et al.* 2019;
31 Prendergast *et al.* 2019; Wadsworth *et al.* 2017) analyses have made a significant contribution to the
32 identification of wild and domestic ungulates remains. Nevertheless, these techniques are not always
33 applicable, as the condition of the faunal remains (*e.g.* poor preservation of DNA and ancient proteins,
34 alteration or modification of the bone surface, lack of reference data) may render them unsuitable for
35 use. Furthermore, they are expensive and can only be used for targeted issues involving a limited
36 number of specimens, not to realize the entire identification of a zooarchaeological collection.
37 Moreover, over the past two decades, geometric morphometrics methods (GMM) have been
38 employed in zooarchaeology to document numerous animal species undergoing domestication
39 processes (Cucchi *et al.*, 2021; Cucchi *et al.* 2023; Evin *et al.* 2015), differentiate between
40 morphologically similar taxa such as sheep and goats (Colominas *et al.* 2019; Haruda 2017; Vuillien
41 2020), and explore species-level variability, such as sheep (Haruda *et al.* 2019; Pöllath *et al.* 2019;
42 Pöllath *et al.* 2019) and deer (Curran 2012). However, GMM are relatively time-consuming, and the
43 observed morphological differences and similarities are based on two- or three-dimensional patterns
44 representing part of the bone being studied and are not a representation of the entire bone.

45

1 The utilisation of machine learning (ML) approaches is being explored for the analysis of biological
2 archives, including bone remains of terrestrial (*e.g.* Moclán *et al.* 2019) and marine mammals (*e.g.*
3 Bickler 2021). Recent studies have evaluated the performance of ML and GMM 2D studies on upper
4 and lower molars of modern and fossil mice (Miele *et al.* 2020; Moclán *et al.* 2023). These studies aim
5 to propose novel identification criteria for these taxa to document the dynamics of human settlements
6 and their role in the emergence and spread of the commensal house mouse (Cucchi *et al.* 2020).
7 Another recent study combined classification methods, including artificial neural networks and GMM
8 2D/3D studies on several teeth and bones of wild and domestic equids and their hybrids (Mohaseb *et al.*
9 *et al.* 2023). The aim was to increase the identification of archaeological equid species and their hybrids
10 in three archaeological sites located in the Middle East. The majority of ML approaches to the
11 identification of ancient animal species rely on the use of two-dimensional images. Nevertheless, the
12 use of three-dimensional imagery in zooarchaeology and, more generally, in archaeology (Andres *et al.*
13 *et al.* 2012; Wyatt-Spratt 2022) has become prevalent, thus enabling the construction and development
14 of osteological and archaeological digital reference collections that ML methods can employ. In this
15 context, the investigation of three-dimensional meshes and point clouds describing biological objects
16 is particularly interesting to mathematicians, both in terms of the object itself and for the complex
17 methodological developments that this represents (Botsch *et al.* 2010; Kazhdan *et al.* 2006; Zhao *et al.*
18 2021). Moreover, the utilisation of 3D imaging, also tried and tested, is pertinent to the issue of species
19 identification, particularly in the context of morphologically similar species such as domestic and wild
20 ruminants.

21
22 This contribution aims to explore the potential of ML approaches directly working on 3D scans and, in
23 particular, on point clouds. Although several ML approaches could serve our purpose (*i.e.* the
24 automatic taxonomic identification from 3D point clouds), some features of the available dataset limit
25 the number of “feasible” approaches. In particular, (i) after the acquisition, the 3D bones have different
26 orientations, scales and number points; (ii) the number of bones for each species is quite limited. Due
27 to the first issue (i), it is difficult to successfully come up with a meaningful notion of distance or
28 similarity between two bones. Several techniques were tested to automatically register (*i.e.* reduce to
29 the same pose and find correspondences between points) collections of point clouds (Evangelidis &
30 Horaud 2017; Myorenko & Song 2010) but, due to the difficulty of the task, such methods failed.
31 Instead, the Iterative Closest Point algorithm (Zhang 2021) would successfully register the whole
32 collection. Still, humans are required to intervene in order to fix a set of benchmarks on each bone
33 manually. This task is long and tedious, and an important aim is to avoid any bias that could be
34 introduced by human intervention at this step. The second feature (ii) prevents from exploiting deep
35 learning architectures (such as (Feng *et al.* 2020; Qi *et al.* 2017), etc.), which need a significant amount
36 of data (here, bones per species) to be appropriately trained. Moreover, deep learning methods would
37 require a downsampling preprocess of the point clouds, whose size is prohibitive for standard
38 architectures, which could lead to information losses.

39
40 In order to fully use the point clouds and highlight morphological features related to the species, as
41 revealed by GMM studies, Topological Data Analysis approach (TDA, Chazal & Michel 2021) was
42 chosen. TDA is a branch of mathematics that studies the structure and the topological properties of
43 data. It has gained popularity in recent years due to its ability to uncover patterns in datasets that are
44 not easily discernible through traditional ML methods (Calsson *et al.* 2008; Dequeant *et al.* 2008;
45 Nicolau *et al.* 2011). The use of TDA as a descriptor for zooarchaeological bones provides a powerful

1 method for addressing the complexities of 3D data analysis. Indeed, as previously mentioned, TDA
2 enables the use of the entire 3D scan, preserving the full point cloud and ensuring that no crucial
3 information is lost. Moreover, the invariance of TDA to isometries such as translation, rotation, and
4 reflection makes it particularly suitable for our study, ensuring that the extracted topological features
5 are intrinsic to the 3D bones and not affected by their positioning.
6

7 Once the topological features of each (3D scan of a) bone are extracted with TDA, each specimen is
8 classified. However, different topological features induce different notions of similarity between
9 bones. Roughly speaking, although two bones are similar in terms of “connected components”, for
10 instance, they might differ in terms of “cycles”. Thus, each topological feature (connected components,
11 cycles, holes) is used to construct a so-called “core” matrix. For instance, the entry (i, j) in the kernel
12 matrix of cycles measures how similar bones i and j are, in terms of cycles. The final objective is to
13 assess the impact of each individual kernel/feature on the classification task (taxonomic identification)
14 and possibly discard features that are redundant. This is precisely what multiple kernel learning (MKL,
15 Gönen & Alpaydin 2011) does. A Bayesian formulation of a logistic regression classifier was used and
16 an original stochastic variational inference approach (SVI, Hoffman *et al.* 2013) was developed to i)
17 perform supervised classification of specimens and ii) evaluate the impact of each topological feature.
18

19 **Research aim**

20
21 The main contribution of this paper is twofold. Firstly, a machine learning routine is proposed to
22 automatically identify morphologically related animal species, like small ruminant herbivores such as
23 *sheep*, *goat*, *Alpine ibex*, *roe deer* and *gazelle*. This routine relies on topological data analysis. Secondly,
24 an original statistical framework is presented, allowing to *weight* the extracted topological features in
25 such a way as to exploit each one of them (*i.e.* multiple kernel learning).

26 Three additional remarks are needed. First, although the classification is *automatic*, the expertise of
27 the zooarchaeologist is required for the analysis of the results to construct an analytical framework
28 and gain an understanding of the way statistics work. Second, the above-mentioned species have been
29 chosen for three main reasons: 1) their morphological proximity; 2) their simultaneous presence in
30 certain archaeological contexts; 3) the large number of taxonomic criteria available in the literature to
31 differentiate some of them, such as sheep and goats, compared to the lack of data for others, such as
32 roe deer and gazelle. Third, the model described here is based on whole bones in a good state of
33 preservation.
34

35 **2. Material and methods**

36 37 **2.1. Astragalus**

38
39 The anatomical part selected for this study is a short bone, the astragalus, from the tarsal joint of the
40 foot (**Supplementary data 1**). It preserves very well in archaeological faunas because it is a small bone,
41 particularly compact and robust and rarely broken intentionally due to its low nutritional value (Barone
42 1976; Popkin *et al.* 2012). These bones presented several anatomical criteria for identifying wild and
43 domestic ungulates discussed by multiple scholars for almost sixty years. The distinction between
44 sheep and goats is well documented (Boessneck *et al.* 1964; Clutton-Brock *et al.* 1990; Fernandez 2001;

1 Prummel & Frish 1986; Salvagno & Albarella 2017; Zeder & Lapham 2010; Zeder & Pilaar 2010), but
 2 still poses problems (Sipilä *et al.* 2023) and represent a challenge as demonstrated by recent GMM
 3 (Gaastra *et al.* 2023; Jeanjean *et al.* 2022; Lloveras *et al.* 2022; Pöllath *et al.* 2019; Vuillien 2020) and
 4 molecular studies (Jeanjean *et al.* 2023; Le Meillour *et al.* 2020). In addition, there are few criteria for
 5 distinguishing between roe deer, gazelle and ibex (Buitenhuis 1988; Crégut-Bonnoure, 2020;
 6 Fernandez 2001; Gudea & Stan 2012; Lavocat 1966; Peters 1989).

7
 8
 9

2.2. 3D models dataset

10 The dataset included 150 3D complete astragali belonging to five taxa: Alpine ibex (*Capra ibex*), sheep
 11 (*Ovis aries*), goat (*Capra hircus*), roe deer (*Capreolus capreolus*) and gazelle (*Gazella cuvieri*, *Gazella*
 12 *dorcas*, *Gazella spekei* and *Gazella sp.*) (**Supplementary data 2**). This dataset does not consider the
 13 specimens' geographical origin or provenance which do not concern our research topic. Gazelle species
 14 are also grouped at the genus level for statistical reasons. The specimens belong to National Museum
 15 of Natural History Mammalian and Birds collection of Paris, CEPAM (UMR 7264) and Archéorient (UMR
 16 5133) labs zooarchaeological reference collection, modern sheep and goat collected for the EvoSheep
 17 collection (ANR-17-CE27-0004), modern goats from BALUT Laboratory Iran, archeological site of
 18 "Grotte de l'Observatoire" Museum of Prehistoric Anthropology of Monaco and archaeological site of
 19 "Tell Sheikh Hassan" Archeorient lab (UMR 5133) (**Table 1**). In order to provide a homogeneous group
 20 for ML analysis, each taxa is represented by 30 specimens (3D astragali).

21

Species	Variety	Curator	Number
Alpine ibex (<i>Capra ibex</i>)	Archaeological from Southern Alps (<i>Liguro-Provençal Bassin</i>)	Museum of Prehistoric Anthropology of Monaco - Archeological site of "Grotte de l'Observatoire"	29
	Modern from Alps	Osteological collection from Thierry Argant (Éveha Lyon, ArAr UMR 5138)	1
Goat (<i>Capra sp.</i>)	Modern <i>Capra nubiana</i> (zoological specimen)	National Museum of Natural History Paris' Mammalian and Birds collection	2
	Modern domestic goat from France	Osteological collections from Archéorient UMR 5133; National Museum of Natural History Paris' Mammalian and Birds collection	5
	Modern domestic goat from Iran	Osteological collections from Bioarchaeology Laboratory, University of Tehran, Iran	20
	Modern domestic goat from Egypt (zoological specimen)	National Museum of Natural History Paris' Mammalian and Birds collection	2
	Modern feral goat from Crete (<i>Capra aegagrus cretica</i>)	Osteological collections from Archéorient UMR 5133	1
Roe deer (<i>Capreolus</i>)	Modern from France	National Museum of Natural History	30

<i>capreolus</i>)		Paris' Mammalian and Birds collection	
Gazelle (<i>Gazella sp.</i>)	Modern <i>Gazella cuvieri</i>	National Museum of Natural History Paris' Mammalian and Birds collection	2
	Modern <i>Gazella dorcas</i>	National Museum of Natural History Paris' Mammalian and Birds collection	5
	Modern <i>Gazella spekei</i>	National Museum of Natural History Paris' Mammalian and Birds collection	1
	Modern <i>Gazella sp.</i>	National Museum of Natural History Paris' Mammalian and Birds collection; Osteological collections from Archéorient UMR 5133	7
	Archaeological from Syria (<i>Gazella cf. subgutturosa?</i>)	Daniel Helmer - Emmanuelle Vila (UMR 5133 Archéorient)	15
Sheep (<i>Ovis aries</i>)	Modern sheep from France	Osteological collections from CEPAM UMR 7264, AASPE UMR 7209 & Archéorient UMR 5133	7
	Modern sheep from Ethiopia	ILRI - Agraw Amane - Emmanuelle Vila - EvoSheep projet (ANR ANR-17-CE27-0004)	23
		Total of 3D astragalus	150

1 Table 1: Summary of sampled modern and archaeological species.

2

3 The astragali were scanned using the Artec Spider blue LED surface scanner and Artec Studio
4 reconstruction software (version 16) and EinScan Pro 2X (Figure 1). The 3D models are reconstructed
5 at a resolution between 0.3 and 0.1 mm using a textured polygonal mesh. Meshes are exported in
6 "ASCII.ply" and "obj." archiving format (Vergniew *et al.* 2017).

7

8 **2.3. Topological data analysis (TDA)**

9

10 Although an in-depth presentation of TDA clearly is outside of the scope of this paper (the interested
11 reader is referred to Chazal & Michel 2021), in this section are sketched the main ideas TDA relies on.

12

13 The input data here is a collection of 150 3D point clouds, and the aim is to extract some useful
14 information (or features) from each point cloud in order to use it to assign the cloud/bone to a species.
15 In order to extract the features, a "continuous" shape is built from the point cloud by progressively
16 connecting data points that are closer to each other with an edge. Here, the reader can safely consider
17 that two points are close if their Euclidean distance in the 3D space is smaller than a given threshold ϵ
18 > 0 . As far as ϵ grows, more and more points are connected and the shape that appears is a collection
19 of simplicial complexes (Figure 2). A simplicial complex can be seen as a higher-dimensional

1 generalization of a neighboring graph: whereas the latter only includes vertices and edges, the former
2 also contains faces, namely triangles and tetrahedrons. The nested family of simplicial complexes that
3 add to each other is called filtration and as long as the family grows, relevant topological features such
4 as connected components, loops and voids are collected via specific methods, such as persistent
5 homology (PH, (Otter *et al.* 2017; Zomorodian & Carlsson 2004) and stored into the so-called
6 persistence diagrams (PDs), that is described in some detail in the next section. It has been shown
7 (Chazal & Michel 2021) that the topological features provide insights into the underlying structure of
8 the data, making TDA particularly useful for analysing highly dimensional and noisy datasets.

9
10 In order to provide the reader with an intuition of what TDA is, the exposition of this pipeline is
11 simplified. However, some remarks are needed:

- 12
13 1. Since the way simplicial complexes are built mainly relies on the relative distance between the
14 points in the cloud, the orientation of the 3D shapes is irrelevant and there is no longer a need
15 to register the collection.
- 16 2. Several notions of distance between points can be chosen and several ways of aggregating
17 simplicial complexes (*i.e.* filtrations) exist.
- 18 3. The topological features are collected in PDs, but alternatives exist (*e.g.* persistence images,
19 barcodes, etc.).

20
21 These few remarks should help the reader to figure out the vastness and richness of topological data
22 analysis.

23 24 **2.3.1. Persistence Diagrams (PDs)**

25
26 In all the experiments, PDs were created from 3D point clouds based on a particular filtration: the
27 Alpha filtration (GUDHI project 2023; Rouvreau 2023). Its main ingredient is a simplicial complex (the
28 Alpha complex) that is built upon the growing balls mechanism that where sketched in the previous
29 section (Figure 2). As said, as the radius ϵ of each ball increases, more simplices fuse and the topology
30 evolves. Persistent homology then tracks these changes, identifying the *birth* and *death* of topological
31 features such as connected components (0-dimensional features), cycles (1-dimensional), and voids
32 (2-dimensional). These “lifelong” are then stored and visualised in a 2D diagram, the PD
33 (Supplementary data 3). Each point's coordinates represent the scale (*i.e.* the value of ϵ) at which a
34 feature appears (birth) and disappears (death) along the filtration process. Three different colors
35 correspond to the different topological features: connected components, in red, loops, in blue and
36 voids, in green. Notice that the value ∞ is allowed on the y-axis and a red point takes this value. It
37 means that at the end of the filtration process a single and huge connected component is still alive,
38 whereas all cycles and voids do not exist anymore.

39
40 PDs were shown to provide a valuable summary of the point cloud's underlying geometrical structure,
41 being robust to perturbation of the data in the Gromov-Hausdorff metric (Chazal & Michel 2021,
42 Theorem 9). The key underlying intuition is that if two bones have very similar shapes, then also their
43 persistence diagrams will be similar and vice versa. Thus, once all the input point clouds are described
44 by their PDs, next step is to compute a pairwise similarity matrix (kernel), whose entry (i, j) is a non-
45 negative number quantifying the similarity between (the PD of) i and (the PD of) j . The kernel matrix

1 can then be used as the input for ML algorithms in order to automatically perform taxonomic
 2 identification (Section 2.5). Next section describes how the similarity between two PDs can be
 3 calculated.

5 **2.4. Discrete optimal transport**

7 To quantify similarities between PDs, several metrics have been developed (Biasotti *et al.* 2011; Efrat
 8 *et al.* 2001). It should be stressed that the notions of similarity and distance between PDs are two sides
 9 of the same coin. Indeed, in general given two data points x and y if the distance $d(x, y)$ between them
 10 is known, a measure of similarity is obtained *via*

$$12 \quad K(x, y) = e^{-d(x,y)/\lambda} \quad (1bis)$$

14 for any real positive λ . That said, the study of distances between discrete probability distributions,
 15 using the Wasserstein distance (Lacombe *et al.*, 2018) from optimal transport has been resorted (OT,
 16 Cuturi 2013). Optimal transport provides an intuitive way to quantify the similarity between two
 17 probability distributions by considering the minimum cost required to transform one distribution into
 18 the other. Here, discrete probability distributions are given consideration. In more details, two
 19 persistent diagrams $P_x = \{x_1 \dots, x_N\}$ and $P_y = \{y_1, \dots, y_M\}$, with N and M denoting the number of
 20 points in each diagram are considered. One can easily define a probability distribution μ_x (respectively
 21 μ_y) by putting mass $1/N$ ($1/M$) over each point in P_x (P_y).

23 **Definition 1:**

25 The Wasserstein distance of order p between μ_x and μ_y is given by

$$26 \quad W_p(\mu_X, \mu_Y) = \left(\min_{\pi \in \Gamma(\mu_X, \mu_Y)} \sum_{i,j} d(x_i, y_j) \pi_{ij} \right)^{1/p}$$

28 where $\Gamma(\mu_x, \nu_x)$ denotes the set of joint probability mass functions (also called couplings or transport
 29 plans) on $P_x \times P_y$ with marginals μ_x and ν_y , respectively, and $d(x, y)$ is the distance between
 30 points x and y in the underlying metric space.

32 Here, $d(x, y)$ are considered to be the Euclidean distance and set $p=2$. Intuitively, one has to imagine
 33 a total mass of 1 is split into N equal portions and distributed over the points of P_x . Now the aim is to
 34 carry all this mass from P_x to P_y under the constraint that, at the end of the day, the total mass is
 35 uniformly distributed over the points of P_y . The Wasserstein distance quantifies the optimal cost of
 36 such an operation, and the optimal transport plan (the one minimising Eq. (1)) tells how much weight
 37 one should move from where to where in order to minimise the effort.

39 However, the PDs gives an additional issue: a green point (for instance) on P_x could be partially or
 40 totally transported to a red or a blue point on P_y , whereas the aim is to keep different topological
 41 features (connected components, loop and voids) well separated. Thus, a PD is “splitted” into three

1 PDs, one for each color (see Figure 4) and compute three Wasserstein distances for each pair of bones,
2 one for each topological feature represented in the PDs.

3
4 Examples of pairwise kernel similarity matrices can be seen in Figures 3 and 5. Kernels were obtained
5 from the Wasserstein distances via Eq. (1), where λ was set equal to the standard deviation of the
6 corresponding distance matrix. The rows and columns of the kernel matrices correspond to the 150
7 bones and the i -th row and j -th column elements denote the Wasserstein similarity between the PDs
8 associated with the i -th and j -th bones, for the corresponding topological features. The main diagonal
9 is the brightest region of each kernel matrix since each bone is trivially at similarity one (and zero
10 distance) from itself. Darker regions in a matrix correspond to higher distances. To narrow down the
11 focus, the 1 and 2-dimensional features (respectively blue and green points in Figure 3) are only
12 considered, which proved to be more informative.

13
14

15 2.5. Supervised multiple kernel learning

16

17 The blueprint of supervised machine learning can be briefly described as follows. Assume a training
18 dataset of $\{x_1, \dots, x_N\}$ observations is given, together with labels $\{y_1, \dots, y_N\}$. Here, the i -th
19 observation x_i is a 3D scan of a bone, in the form of a point cloud, and its label y_i can be seen as an
20 integer ranging from 1 to Q and labelling the species of the bone. As seen in Section 2.2, $N = 150$ and
21 $Q = 5$ for us. Next ingredient is a function of the data (the classifier), say f_θ , depending on some
22 parameters θ and associating to each observation x_i a predicted label $f_\theta(x_i) =: \hat{y}_i$, *i.e.* another integer
23 between 1 and Q . then f_θ is “trained” by solving the following minimisation problem

24

$$25 \min_{\theta} \left(\sum_{i=1}^N L(y_i, \hat{y}_i) \right).$$

26

27 where $L(\cdot, \cdot)$ is a loss function. So, roughly speaking, the value of θ should be such that the mismatch
28 between the predicted and the actual labels is minimal, on average. Once θ is optimised, the final aim is
29 to be able to correctly predict the label y^* of a new test data point x^* , via $f_\theta(x^*)$. In order to check that it
30 is actually the case, in these experiments, all dataset ($N = 150$) has been split into train ($N_{\text{train}} = 120$) and
31 test ($N_{\text{test}} = 30$) and the accuracy (*i.e.* the proportion of correctly identified specimens) reported on the
32 test dataset. To assess the robustness of the test accuracy, 50 random train/test splits are performed
33 allowing the computation of a mean accuracy and a standard deviation.

34

35 Now, one way to define f_θ is to pass through kernel matrices. Several ML classifiers are specifically
36 designed to leverage the representation of the data through kernel matrices. Popular methods include
37 support vector machines (SVM), kernel discriminant analysis (KDA) and kernel logistic regression (KLR).
38 An in-depth description of such methods is outside the scope of this paper and the interested reader is
39 referred to (Hastie *et al.* 2009, Chapters 4,5 and 12).

40

41 In this work, the focus is on multi-class KLR with multiple kernels. Indeed, as mentioned above,
42 different types of filtrations and/or topological features lead to different similarity matrices, each of
43 which could contain different discriminant information and relying on a single kernel matrix chosen by

1 the user could be not the optimal strategy. Multiple kernel learning (MKL) approaches [47] are
2 specifically designed to manage this kind of situation: they allow one to properly weigh the input
3 kernels and possibly discard the useless ones. Formally, D different kernel matrices $\{K(1), \dots, K(D)\}$ with
4 $K(d) \in \mathbb{R}^{N \times N}$ were considered for each d . The aim is for an optimal convex combination of such
5 matrices, defined by

$$6 \quad K := \sum_{d=1}^D \beta_d K^{(d)} \in \mathbb{R}^{N \times N} \text{ with } \beta_d \geq 0 \text{ for each } d. \quad (2)$$

7
8
9 By adopting a hybrid Bayesian formulation, the kernel weights β_d are treated as random variables,
10 following an Exponential prior distribution, whereas the other weights intervening in the KLR are
11 treated as parameters to optimise. This approach allows to revisit MKL in an original way, being a com-
12 promise between pure optimization strategies (Rakotomamonjy *et al.* 2008) and fully Bayesian ones
13 (Damoulas & Girolami 2008; Gonen 2012). Moreover, recent developments in importance-weighted
14 stochastic variational inference (Sobolev & Vetrov 2019) have enabled the optimization problem in a
15 fully differentiable manner, via Stochastic Gradient Descent (SGD, Bottou 2010), and perform posterior
16 inference on β_1, \dots, β_D .

17 18 **3. Results**

19 20 **3.1 Distance matrices (TDA) and optimal transport distances**

21
22 The topological dimension 1 (loops) without normalisation indicates that Alpine ibex, roe deer, and
23 gazelle respectively form dense clusters, in contrast to sheep and goat (Figure 3a). Sheep and goat
24 exhibit an intra-specific significant topological heterogeneity, revealed by a marked color gradient from
25 dark blue to yellow. This is mainly due to both the breed factor and the geographical heterogeneity of
26 the specimens studied. For sheep, specimens from Ethiopian breeds (in dark green and yellow) are
27 distinct from specimens from French breeds (in dark blue) (Figure 4). The same is true for goat.
28 Specimens attributed to breeds of Iranian origin (in dark blue) differ from specimens of French breeds
29 (in green and yellow) except for one specimen (number 57) (Supplementary data 4). Furthermore,
30 although ibex form a homogeneous cluster, one individual stands out (dark blue): this is the only
31 modern ibex in the dataset. When comparing species, the same topological dimension reveals a high
32 degree of separation between wild and domestic species. In more details: ibex form a single
33 community; roe deer look quite similar to gazelle (two wild taxa) and sheep often look indistinguishable
34 from goat (two domestic species) despite no apparent link to their breed or geographical origin.
35 Nevertheless, there is a certain topological proximity between roe deer and gazelle, on one side, and
36 specimens of French sheep and goat, on the other (whereas the morphology of gazelle and roe deer is
37 clearly distinguishable from that of African sheep and Asian goat).

38
39 Also, the topological dimension 2 (holes) without normalisation highlights differences between Alpine
40 ibexes and other species (Figure 3b). These differences could be explained by the size of the astragalus.
41 Ibex are taller than the other species in the dataset. However, if the size effect were crucial, gazelle,
42 the smallest taxon in the dataset, should be clearly distinguishable from the other taxa. This is not the
43 case. Consequently, although this parameter is undoubtedly important, it does not appear to be the
44 sole factor responsible for the structuring of the dataset. Interestingly, whereas in the previous matrix
45 roe deer and gazelle looked quite similar, here the difference between them is more accentuated

(although an increased similarity with goat appears). This point clearly illustrates why the adoption of several kernel matrices is beneficial to the taxonomic identification at the species level. The topological dimension 1 (loops) with bone normalisation allows for the clear distinction of Alpine ibex, roe deer and gazelle from groups of domestic caprine (sheep and goat) (Figure 5a). This outcome indicates that normalisation does not directly impact classification at the inter-specific level.

Finally, the topological dimension 2 (holes) with bone normalisation might seem of no particular interest (Figure 5b) since the similarities between specimens of each species render it impossible to distinguish between the species in question (except for roe deer, partially). However, as demonstrated in the next section, supervised MKL remarkably benefits from this matrix since it is the only one highlighting similarities between specimens of the same species far from each other in previous representations, especially sheep and goat.

3.2 Classification: Supervised MKL and automatic taxonomic identification of bones

For the supervised MKL part, as mentioned, the whole dataset (N = 150 bones) was randomly divided into 120 (80%) as train dataset and 30 (20%) as test dataset. It is recalled that such a random split is repeated 50 times. Table 2 reports the average test accuracy together with its standard deviation (in small characters). The second column reports the global test accuracy, whereas the remaining columns show the test accuracies for each species.

	Acc. _{std.dv.}	Alpine ibex	Goat	Roe deer	Gazelle	Sheep
Multiple KLR	0.811 _{0.064}	0.973 _{0.061}	0.623 _{0.218}	0.927 _{0.101}	0.897 _{0.124}	0.637 _{0.228}

Table 2: Average test classification accuracy.

In order to better inspect how the test accuracy behaves as a function of the train/test data split, a Box-and-Whisker plot is provided in Figure 6a. From the bottom to the top, each Box-and-Whisker reports for each species the minimum test accuracy, the lower quartile, the median accuracy, the upper quartile, and the maximum test accuracy. Outliers are represented as single points. The only modern specimen present in the Alpine ibex dataset is the outlier. The gazelle outlier corresponds to two specimens: a modern specimen (specimen number 1992 1844) and an archaeological specimen from Tell Sheikh Hassan (specimen number TSH 7Z23).

The average test accuracy of 81,1% (sensibly higher than 15% that one would obtain from a random assignment) is boosted by the accuracy on the wild species, which our approach can successfully identify. Figure 6a shows that the median test accuracy on wild species is 100%. The mean/median accuracies for sheep and goat are sensibly lower and come with a higher variability. However, as pointed out when describing the kernel matrices in Section 3.1, the non-homogeneous intra-specific blocks for these two species correspond to different breeds.

Table 3 shows the weights β_1, \dots, β_4 , introduced in Eq. (2) and estimated via importance-weighted stochastic variational inference (IW-SVI). In more detail, our Bayesian inference strategy allows us to sample from the approximate posterior distribution of β_1, \dots, β_4 , given the data and the other model

1 parameters. Table 3 reports the posterior mean of each weight with its standard deviation (in small
 2 character). Kernel density estimates of weights' posterior distributions can be seen in Figure 6b. The
 3 mode of each estimated density (*i.e.* the maximum a posteriori estimate of the corresponding weight)
 4 is away from zero, meaning that multiple KLR exploits all the topological features considered, either
 5 for normalized or unnormalized point clouds. Interestingly, the most important weight is put on β_4 ,
 6 corresponding to the apparently less expressive kernel matrix of Figure 8b. This happens precisely
 7 because this kernel matrix is the only one allowing the algorithm to densify the clusters of sheep and
 8 goat, respectively. Furthermore, the weight given to each matrix indicates the influence of size on
 9 classification. If size was the more discriminating factor, then β_1 (non-standardised matrix) should have
 10 a greater weight than β_3 (standardised matrix). The results show the opposite (Table 3), indicating that
 11 bone size is not the more discriminating factor for the final classification.

	Kernel weights (posterior mean)			
	β_1	β_2	β_3	β_4
IW-SWI	0.231 _{0.007}	0.233 _{0.006}	0.275 _{0.010}	0.481 _{0.020}

17 Table 3: Kernel weight's estimates.

18 This can clearly be seen in Figure 7 which shows the learned kernel matrix K introduced in Eq. (2) and
 19 obtained as a linear combination of the Wasserstein kernels, weighted via the optimal betas shown in
 20 Table 3. K exhibits brighter diagonal blocks (especially for wild species) and the blocks corresponding
 21 to sheep and goats look more similar to dense blocks with respect to what Figures 3 and 5 show.

23 4. Discussion

24 4.1. Inter-specific identification: wild VS domestic

25 About the main research question of the inter-specific identification of wild and domestic small
 26 ruminants, TDA proved to be very proficient at correctly identifying wild specimens from a 3D scan of
 27 their astragalus, whereas it suffers with the identification of domestic specimens.

28 In other words, TDA makes it possible to clearly distinguish ibex, roe deer, and gazelle astragalus from
 29 sheep and goat astragalus, which underlines its value. Indeed, it is often difficult to distinguish these
 30 wild animals from sheep and goats in archaeological contexts where all these species are represented
 31 and where anatomical criteria are not discriminating. In contrast, the use of TDA does not address the
 32 challenges faced by the zooarchaeological community in distinguishing sheep from goats. Here, the
 33 result is dependent on the discriminative capacity of TDA and the dataset; the obtained classification
 34 exhibits the intraspecific morphological variability of the selected sheep and goat breeds.

35 Nevertheless, the fact that the median accuracy is higher than the mean for the wild species and the
 36 presence of outliers in Alpine ibex and gazelle is explained by the heterogeneity of our dataset: a high
 37 intra-specific variability of a few specimens lowers the mean accuracy. For instance, the single modern
 38 Alpine ibex in the collection turns out to be dissimilar from the other ibexes. For example, the unique

1 modern Alpine ibex that we have in our collection. As we saw in the figures in Section 3, it is dissimilar
2 from the other ibexes; with regard to the train dataset, multiple KLR can safely identify all the
3 archaeological ibex in the test. Conversely, when the modern specimen is included in the test dataset,
4 the classifier fails to identify it (since trained on archaeological ibex dated to the Upper Pleistocene),
5 and turns it into an outlier in the Alpine ibex group. This outcome is consistent with the difficulties
6 encountered by zooarchaeologists in using modern datasets to identify ancient animal populations.
7 This is due to the contrast between the number of current species and subspecies compared to fossil
8 species and their morphological diversity (Crégut-Bonnoure 2020; Crégut-Bonnoure & Fernandez
9 2018; Urena *et al.* 2018). However, our findings indicate that this approach should be employed with
10 great caution when attempting to address the question of the evolution of morphology. Indeed, upon
11 noting the presence of two outliers within the gazelle group, a more detailed examination reveals that
12 interpreting these differences in relation to the available temporal origin information is challenging.
13 The dissimilarities between modern and archaeological gazelles (Holocene, between 9000 and 6000
14 BC) are not immediately apparent: TDA makes it difficult to identify variability among the gazelle
15 species. This is likely due to the dataset itself, which does not accurately reflect the morphological
16 variability among the modern species (*Gazella cuvieri*, *Gazella dorcas*, *Gazella spekei*). However, we
17 intend to compare available proteomic data (Culley *et al.* 2021; Janzen *et al.* 2021; Le Meillour *et al.*
18 2020; Le Meillour *et al.* 2023) and future morphological and morphometric studies (Vuillien 2024).

19 **4.2. Beyond the classification problem: is astragalus a good taxonomic marker?**

20 The results obtained for species classification prompt the question of whether the astragalus can be
21 employed as an interspecific identification marker. The geographical provenance of domestic and wild
22 specimens appears to exert a more pronounced impact than the distinction between species. As we
23 demonstrated, TDA cannot see sheep and goats as uniform and separated clusters (within the limits
24 of the explored filtrations) and this fact has some significant consequences for the weighting of the
25 kernels. The misidentification of sheep and goat species is attributed to the morphological variability
26 of domestic breeds and their geographical origin. In addition, the difference in the ibex group observed
27 between the single modern specimen and the other archaeological specimens may be correlated with
28 the morphological evolution of the Alpine ibex over time and would point to the potential of the
29 astragalus as an ecomorphological marker (Barr 2014; DeGusta 2003; Plummer *et al.* 2008). The
30 biometrical and biomolecular studies carried out for this species show its morphotypic diversity over
31 time, linked to environmental changes, the rocky areas frequented and human pressure (Crégut-
32 Bonnoure 2020). However, it is still difficult to identify these factors precisely from faunal remains.
33 TDA could be a valuable tool for exploring morphological variability at the intra-specific level for
34 domestic and wild species and their adaptation across time and the environment, such as recent GMM
35 studies applied to the same modern sheep breeds (Bader *et al.* 2022) and archaeological Alpine ibex
36 populations. In archaeozoology, identifying and classifying of domestic sheep and goat morphotypes
37 is of great importance. Such an analysis can provide insight into the evolution of zootechnical practices,
38 economic development, and human society (Vila *et al.* 2021).

39 **4.3. More robust assessment of the result obtained via TDA**

40 To correctly estimate the potential of the method proposed in this paper, it will be crucial in the future
41 to analyse the same dataset with other approaches either directly based on human expertise to have
42 a “human” accuracy relying on anatomical criteria or automatic, such as GMMs. This paper aimed to
43 test TDA and supervised MKL on a zooarchaeological dataset, nevertheless, TDA features would

1 certainly express the best of their potential in conjunction with other features, such as anatomical
2 criteria and/or GMM patterns.

3 Another crucial aspect of machine learning routines in general is *explainability*: to know whether a
4 classifier can be trusted, we need to understand how it works (on this topic, see Rudin 2019).
5 Unfortunately, whereas TDA keeps track of the lifelong of topological features such as connected
6 components, loops, and holes, it is currently not possible to know where these features are located on
7 the surface of a 3D bone. This makes it very difficult, for instance, to assess whether the similarities that
8 our approach detected between French sheep and roe deer/gazelle are based on meaningful and
9 previously unexplored morphological patterns or if they are due to some other geometric components.
10 Although some model agnostic explaining techniques exist (Lundberg & Lee 2017; Ribeiro *et al.* 2016),
11 their use in the context of 3D point clouds, in conjunction with TDA, is not immediate at all. This issue
12 will be addressed in future research to propose new anatomical features. Finally, it will be of interest
13 to test deep learning models. Although requiring a huge amount of training data, such models can obtain
14 impressive results, and their behaviors could be more easily captured than the one of TDA, either via
15 explaining techniques or ad-hoc architectures.

16 **Conclusion and perspectives**

17 This paper mainly focuses on the taxonomic identification of wild and domestic small ruminants from
18 3D scans of complete astragalus of modern and archaeological specimens. The problem was framed
19 as a supervised learning problem and addressed using TDA, optimal transport and an original inference
20 routine. From one side, the topological features extracted with TDA proved to be very discriminant in
21 classifying wild species (Alpine ibex, roe deer, gazelle). The main strengths of the proposed approach
22 are a median test accuracy of 100% for these species and the fact that our routine is entirely automated
23 (the expert's intervention is only required to analyse the results). On the other hand, TDA/MKL partly
24 failed to identify the modern domestic species goat and sheep. However, an in-depth analysis of the
25 reasons for such difficulty revealed that TDA might be better suited for the intra-specific classification
26 of such species, for which our method seems likely to perceive a lot of detail. Another drawback of
27 TDA methods is their lack of explainability: it is not possible to know which part of the bone contributed
28 the most to the identification. In light of the above remarks, a few avenues for future research can be
29 outlined: (i) the creation of ad-hoc datasets of 3D scans to assess the capabilities of TDA/MKL intra-
30 specific classifiers (ii) combine TDA features to anatomical criteria and GMM morphological patterns;
31 (iii) test deep learning methods on an increasing dataset to test others classification approaches. In
32 conclusion, this research demonstrates the effectiveness of the TDA/MKL methods in extracting
33 authentic biological information that can be interpreted by archaeozoologists, as well as novel
34 information that cannot be detected by traditional anatomical criteria. In this sense, these approaches
35 represent a milestone in the dialogue between two scientific disciplines, mathematics and
36 archaeology, by providing information comparable to that obtained by palaeogenetics.

37 **Acknowledgements**

38
39 This paper constitutes the doctoral thesis in applied mathematics of the co-first author and as such
40 this project has received financial support from the CNRS through the MITI interdisciplinary programs.
41 It also forms a part of the postdoctoral fellowship "IBEX" supported by the French government through
42 the France 2030 investment plan managed by the National Research Agency (ANR), as part of the

1 Initiative of Excellence of Université Côte d'Azur under reference number ANR-15-IDEX-01. We would
2 like to thank Arch'Al'Story project (Ministère de l'Enseignement Supérieur et de la Recherche and
3 University Côte d'Azur) for funding this project. We would like to extend our gratitude to the
4 researchers from the Museum of Prehistoric Anthropology in Monaco, the Laboratoire Départemental
5 de Préhistoire du Lazaret in Nice (France), the National Museum of Natural History in Paris (France)
6 and BioArch (UMR 7209, Paris) and Archéorient (UMR 5133, Lyon) research laboratories for their
7 invaluable contributions. We are grateful to the ANR EvoSheep project (ANR-17-CE27-0004) for
8 providing access to the modern sheep and goat collection collected in 3D during the first postdoctoral
9 position of the first author to use part of it in this study. We are grateful to the collective research
10 project "Paleoecology of the Lazaret Cave: human-environment interactions on the coast of the
11 meridional Alps during the late Middle Pleistocene (MIS 6)," granted by the DRAC PACA (French
12 Ministry of Culture) that supported work in Pleistocene archaeological collection and provided the
13 access to the 3D Einscan. We would like to express our gratitude to Emmanuel Desclaux for his support
14 during the 3D modelling of Alpine ibex collection. We also express our gratitude to the "Imagery
15 platform in bioarchaeology" coordinated by Thomas Cucchi at the BioArch laboratory. We would also
16 like to thank Joséphine Lesur for her assistance and authorisation to provide samples of the modern
17 gazelle and roe deer collection held at the National Museum of Natural History of Paris.

18

19 **Declaration of interest:** none

20

21 **References**

22

23 Adcock, A., Carlsson, E. and Carlsson, G. 2013. The ring of algebraic functions on persistence bar codes.
24 *arXiv preprint arXiv:1304.0530*.

25

26 Alberto, F.J., Boyer, F., Orozco-terWengel, P., Streeter, I., Servin, B., De Villemeureuil, P., Benjelloun, B.,
27 Librado, P., Biscarini, F., Colli, .L, Barbato M., Zamani W., Alberti A., Engelen S., Stella A., Joost S.,
28 Ajmone-Marsan P., Negrini R., Orlando L., Reza Rezaei H., Naderi S., Clarke L., Flicek P., Wincker P.,
29 Coissac E., Kijas J., Tosser-Klopp G., Chikhi A., Bruford W. M., Taberlet P., Pompanon F. 2018.
30 Convergent genomic signatures of domestication in sheep and goats. *Nature communications* 9(1):
31 813. DOI: 10.1038/s41467-018-03206-y

32

33 Andrés, AN., Pozuelo, FB., Marimón, JR. and de Mesa Gisbert, A. 2012. Generation of virtual models of
34 cultural heritage. *Journal of Cultural Heritage* 13(1): 103–106.

35

36 Bader, C., Mallet, C., Chahoud, J., Amane, A., De Cupere, B., Berthon, R, Lavenne, F., Mohaseb, A.,
37 Davoudi, H., Albesso, M., Fathi H., Vuillien M., Lesur J. Helmer D., Gourichon L., Hanotte O., Mashkour
38 M., Vila E., Cucchi T. 2022. Are petrous bones just a repository of ancient biomolecules? Investigating
39 biosystematic signals in sheep petrous bones using 3D geometric morphometrics. *Journal of*
40 *Archaeological Science: Reports* 43: 103447. DOI : 10.1016/j.jasrep.2022.103447.

41

42 Barone, R. 1976. *Anatomie comparée des mammifères domestiques*. Ostéologie deuxième édition
43 revue et augmentée. Vigot. Paris.

44

1 Barr, WA. 2014. Functional morphology of the bovid astragalus in relation to habitat: Controlling
2 phylogenetic signal in ecomorphology. *Journal of Morphology* 275(11): 1201–1216. DOI:
3 <https://doi.org/10.1002/jmor.20279>.

4

5 Biasotti, S., Cerri, A., Frosini, P. and Giorgi, D. 2011. A new algorithm for computing the 2-dimensional
6 matching distance between size functions. *Pattern Recognition Letters* 32(14): 1735–1746.

7

8 Bickler, SH. 2021. Machine Learning Arrives in Archaeology. *Advances in Archaeological Practice* 9(2):
9 186–191. DOI: <https://doi.org/10.1017/aap.2021.6>.

10

11 Boessneck, J., Müller, H-H. and Teichert, M. 1964. *Osteologische Unterscheidungsmerkmale zwischen*
12 *Schaf (Ovis aries Linné) und Ziege (Capra hircus Linné)*. Verlag nicht ermittelbar.

13

14 Bonneel, N., Rabin, J., Peyré, G. and Pfister, H. 2015. Sliced and radon wasserstein barycenters of
15 measures. *Journal of Mathematical Imaging and Vision* 51: 22–45.

16

17 Botsch, M., Kobbelt, L., Pauly, M., Alliez, P. and Lévy, B.. 2010. *Polygon mesh processing*. CRC press.

18

19 Bottou, L. 2010. Large-scale machine learning with stochastic gradient descent. In: *Proceedings of*
20 *COMPSTAT'2010: 19th International Conference on Computational Statistics Paris France, August 22-*
21 *27, 2010 Keynote, Invited and Contributed Papers*. 2010. Springer. pp. 177–186.

22

23 Buitenhuis, H. 1988. *Archeozoölogisch onderzoek langs de Midden-Eufrat: onderzoek van het*
24 *faunamateriaal uit zes nederzettingen in Zuidoost-Turkije en Noord-Syrië daterend van ca. 10.000 BP*
25 *tot 1400 AD*.

26

27 Carlsson, G., Ishkhanov, T., De Silva, V. and Zomorodian, A. 2008. On the local behavior of spaces of
28 natural images. *International journal of computer vision* 76: 1–12.

29

30 Chazal, F. and Michel, B. 2021. An introduction to topological data analysis: fundamental and practical
31 aspects for data scientists. *Frontiers in artificial intelligence* 4: 667963.

32

33 Clutton-Brock, J., Dennis-Bryan, K., Armitage, PL. and Jewell, PA. 1990. Osteology of the Soay sheep.
34 *Bulletin british Museum natural History Zoology* 1(56): 1–56.

35

36 Colominas, L., Evin, A., Burch, J., Campmajó, P., Casas, J., Castanyer, P., Carreras, C., Guardia, J., Olesti,
37 O., Pons, E., Tremoleda, J. and Palet, J-M. 2019. Behind the steps of ancient sheep mobility in Iberia:
38 new insights from a geometric morphometric approach. *Archaeological and Anthropological Sciences*
39 11(9): 4971–4982. DOI: <https://doi.org/10.1007/s12520-019-00837-0>.

40

41 Crégut-Bonnoure, E. 2020. *Les Ovibovini, Caprini et Ovini (Mammalia, Artiodactyla, Bovidae, Caprinae)*
42 *du Plio-Pléistocène d'Europe*. BAR Publishing International Series. BAR International Series.

43

44 Crégut-Bonnoure, E. and Fernandez, P. 2018. Perspectives morphométriques et phylogéniques du
45 genre *Capra* au Pléistocène (Mammalia, Artiodactyla, Caprinae). *Quaternaire* 29(3): 243–254.

46

47 Cucchi, T., Domont, A., Harbers, H., Evin, A., Alcàntara Fors, R., Saña, M., Leduc, C., Guidez, A., Bridault,
48 A., Hongo, H., Price M., Peters J., Briois F., Guilaine J. Vigne J.-D. 2021. Bones geometric morphometrics

1 illustrate 10th millennium cal. BP domestication of autochthonous Cypriot wild boar (*Sus scrofa circeus*
2 *nov. ssp.*). *Scientific Reports* 11(1): 11435. DOI: 10.1038/s41598-021-90933-w
3
4 Cucchi, T., Harbers, H., Neaux, D., Balasse, M., Garbé, L., Fiorillo, D., Bocherens, H., Drucker, D., Zanolli,
5 C., Cornette, R., Arbogast R.-M., Bréhard S., Bridault A. Gourichon L., Guilaine J., Manen C., Perrin T.,
6 Schafberg R., Tresset A., Vigne J.-D., Herrel A. 2023. 4500 years of morphological diversification in
7 Western Europe wild boars (*Sus scrofa*) and the consequences of the Neolithic transition. *Quaternary*
8 *Science Reviews* 309: 108100. DOI: 10.1016/j.quascirev.2023.108100
9
10 Cucchi, T., Papayianni, K., Cersey, S., Aznar-Cormano, L., Zazzo, A., Debruyne, R, Berthon, R., Bălăşescu,
11 A., Simmons, A., Valla, F., Hamilakis Y., Mavridis F. Mashkour M., Darvish J., Siahsarvi R., Biglari F.,
12 Petrie A. C., Weeks L., Sardari A., Maziar S., Denys C., Orton D., Jenkins E., Zeder M., Searle J. B., Larson
13 G., Bonhomme F., Auffray J.-C., Vigne J.-D. 2020. Tracking the Near Eastern origins and European
14 dispersal of the western house mouse. *Scientific reports* 10(1): 8276. DOI: 10.1038/s41598-020-64939-
15 9
16
17 Culley, C., Janzen, A., Brown, S., Prendergast, M.E., Shipton, C., Ndiema, E, Petraglia, M.D., Boivin, N.
18 and Crowther, A. 2021. Iron Age hunting and herding in coastal eastern Africa: ZooMS identification of
19 domesticates and wild bovids at Panga ya Saidi, Kenya. *Journal of Archaeological Science* 130: 105368.
20 DOI: <https://doi.org/10.1016/j.jas.2021.105368>.
21
22 Curran, SC. 2012. Expanding ecomorphological methods: geometric morphometric analysis of Cervidae
23 post-crania. *Journal of Archaeological Science* 39(4): 1172–1182.
24
25 Cuturi, M. 2013. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural*
26 *information processing systems* 26, NIPS'13: Proceedings of the 27th International Conference on
27 Neural Information Processing Systems: 2292-2300.
28
29 Daly, K.G., Maisano Delsler, P., Mullin, V.E., Scheu, A., Mattiangeli, V., Teasdale, M.D., Hare, A.J., Burger,
30 J., Verdugo, M.P., Collins, M.J., Kehati, R., Erek, C.M., Bar-Oz, G., Pompanon, F., Cumer, T., Çakırlar, C.,
31 Mohaseb, A.F., Decruyenaere, D., Davoudi, H., Çevik, Ö., Rollefson, G., Vigne, J.-D., Khazaeli, R., Fathi,
32 H., Doost, S.B., Rahimi Sorkhani, R., Vahdati, A.A., Sauer, E.W., Azizi Kharanaghi, H., Maziar, S.,
33 Gasparian, B., Pinhasi, R., Martin, L., Orton, D., Arbuckle, B.S., Benecke, N., Manica, A., Horwitz, L.K.,
34 Mashkour, M. and Bradley, D.G. 2018. Ancient goat genomes reveal mosaic domestication in the
35 Fertile Crescent. *Science* 361(6397): 85–88. DOI: <https://doi.org/10.1126/science.aas9411>.
36
37 Damoulas, T. and Girolami, MA. 2008. Probabilistic multi-class multi-kernel learning: on protein fold
38 recognition and remote homology detection. *Bioinformatics* 24(10): 1264–1270.
39
40 DeGusta, D. and Vrba, E. 2003. A method for inferring paleohabitats from the functional morphology
41 of bovid astragali. *Journal of Archaeological Science* 30(8): 1009–1022. DOI:
42 [https://doi.org/10.1016/S0305-4403\(02\)00286-8](https://doi.org/10.1016/S0305-4403(02)00286-8).
43
44 Dequeant, M.-L., Ahnert, S., Edelsbrunner, H., Fink, T.M., Glynn, E.F., Hattem, G., Kudlicki, A., Mileyko,
45 Y., Morton, J., Mushegian, A.R., Lior P. Rowicka M., Shiu A., Sturmfels B., Pourquié O. 2008. Comparison
46 of pattern detection methods in microarray time series of the segmentation clock. *PLoS One* 3(8):
47 e2856. DOI: [10.1371/journal.pone.0002856](https://doi.org/10.1371/journal.pone.0002856)
48

1 Efrat, A., Itai, A. and Katz, M.J. 2001. Geometry helps in bottleneck matching and related problems.
2 *Algorithmica* 31: 1–28.
3

4 Evangelidis, G.D. and Horaud, R. 2017. Joint alignment of multiple point sets with batch and
5 incremental expectation-maximization. *IEEE transactions on pattern analysis and machine intelligence*
6 40(6): 1397–1410.
7

8 Evin, A., Girdland Flink, L., Balasescu, A., Popovici, D., Andreescu, R., Bailey, D., Mirea, P., Lazar, C.,
9 Boroneant, A., Bonsall, C., Strand Vidarsdottir, U., Brehard, S., Tresset, A., Cucchi, T., Greger, L. and
10 Dobney, K.. 2014. Unravelling the complexity of domestication: a case study using morphometrics and
11 ancient DNA analyses of archaeological pigs from Romania. *Philosophical Transactions B* 370(1660): 1–
12 8. DOI: <https://doi.org/10.1098/rsb.2013.0616>.
13

14 Fabrizi, I., Flament, S., Delhon, C., Gourichon, L., Vuillien, M., Oueslati, T., Auguste, P., Rolando, C. and
15 Bray, F. 2024. Low-Invasive Sampling Method with Tape-Disc Sampling for the Taxonomic Identification
16 of Archeological and Paleontological Bones by Proteomics. *Journal of Proteome Research* acs.
17 jproteome.4c00083. DOI: <https://doi.org/10.1021/acs.jproteome.4c00083>.
18

19 Feng, M., Zhang, L., Lin, X., Gilani, SZ and Mian, A. 2020. Point attention network for semantic
20 segmentation of 3D point clouds. *Pattern Recognition* 107: 107446.
21

22 Fernandez, H. 2001. *Ostéologie comparée des petits ruminants eurasiatiques sauvages et domestiques*
23 *(genres Rupicapra, Ovis, Capra et Capreolus): diagnose différentielle du squelette appendiculaire*.
24 Thèse de Doctorat. Muséum d’histoire naturelle Genève, Université de Genève.
25

26 Flamary, R., Courty, N., Gramfort, A., Alaya, M.Z., Boisbunon, A., Chambon, S., Chapel, L., Corenflos, A.,
27 Fatras, K., Fournier, N., Gautheron, L., Gayraud, N.T.H., Janati, H., Rakotomamonjy, A., Redko, I., Rolet,
28 A., Schutz, A., Seguy, V., Sutherland, D.J., Tavenard, R., Tong, A. and Vayer, T. 2021. POT: Python
29 Optimal Transport. *Journal of Machine Learning Research* 22(78): 1–8.
30

31 Gaastra, JS. 2023. Corrigendum to “Domesticating details: 3D geometric morphometrics for the
32 zooarchaeological discrimination of wild, domestic and proto-domestic sheep (*Ovis aries*) and goat
33 (*Capra hircus*) populations” [J. Archaeol. Sci. 151 (2023) 105723]. *Journal of Archaeological Science*
34 153: 105768. DOI: <https://doi.org/10.1016/j.jas.2023.105768>.
35

36 Gönen, M. 2012. Bayesian efficient multiple kernel learning. *arXiv preprint arXiv:1206.6465*.
37

38 Gönen, M. and Alpaydın, E. 2011. Multiple kernel learning algorithms. *The Journal of Machine Learning*
39 *Research* 12: 2211–2268.
40

41 Gudea, A. and Stan, F. 2012. The discriminative macroscopical identification of the bones of sheep
42 (*Ovis aries*), goat (*Capra hircus*) and roe deer (*Capreollus capreollus*). 2. Elements of the hindlimb
43 skeleton. *Bulletin UASMV-Ve Med* 69: 132–139.
44

45 Haruda, AF. 2017. Separating sheep (*Ovis aries* L.) and goats (*Capra hircus* L.) using geometric
46 morphometric methods: an investigation of Astragalus morphology from late and final Bronze age
47 central asian contexts. *International Journal of Osteoarchaeology* 27(4): 551–562.
48

1 Haruda, A.F., Varfolomeev, V., Goriachev, A., Yermolayeva, A. and Outram, A.K. 2019. A new
2 zooarchaeological application for geometric morphometric methods: Distinguishing *Ovis aries*
3 morphotypes to address connectivity and mobility of prehistoric Central Asian pastoralists. *Journal of*
4 *Archaeological Science* 107: 50–57. DOI: <https://doi.org/10.1016/j.jas.2019.05.002>.
5
6 Hastie, T., Tibshirani, R., Friedman, J.H. and Friedman, J.H. 2009. *The elements of statistical learning:*
7 *data mining, inference, and prediction*. Springer.
8
9 Hoffman, M.D., Blei, D.M., Wang, C. and Paisley, J. 2013. Stochastic variational inference. *Journal of*
10 *Machine Learning Research*.
11
12 Janzen, A., Richter, K.K., Mwebi, O., Brown, S., Onduso, V., Gatwiri, F., Ndiema, E., Katongo, M.,
13 Goldstein, S.T., Douka, K. and Boivin, N. 2021. Distinguishing African bovids using Zooarchaeology by
14 Mass Spectrometry (ZooMS): New peptide markers and insights into Iron Age economies in Zambia
15 Adams, JW (ed.). *PLOS ONE* 16(5): e0251061. DOI: <https://doi.org/10.1371/journal.pone.0251061>.
16
17 Jeanjean, M., Haruda, A., Salvagno, L., Schafberg, R, Valenzuela-Lamas, S., Nieto-Espinet, A., Forest,
18 V., Blaise, E., Vuillien, M., Mureau, C. and Evin, A. 2022. Sorting the flock: Quantitative identification of
19 sheep and goat from isolated third lower molars and mandibles through geometric morphometrics.
20 *Journal of Archaeological Science* 141: 105580. DOI: <https://doi.org/10.1016/j.jas.2022.105580>.
21
22 Jeanjean, M., McGrath, K., Valenzuela-Lamas, S., Nieto-Espinet, A., Schafberg, R., Parés-Casanova,
23 P.M., Jiménez-Manchón, S., Guintard, C., Tekkouk, F., Ridouh, R., Mureau, C. and Evin, A. 2023. ZooMS
24 confirms geometric morphometrics species identification of ancient sheep and goat. *Royal Society*
25 *Open Science* 10(9): 230672. DOI: <https://doi.org/10.1098/rsos.230672>.
26
27 Kazhdan, M., Bolitho, M. and Hoppe, H. 2006. Poisson surface reconstruction. In: *Proceedings of the*
28 *fourth Eurographics symposium on Geometry processing*. 2006. p.
29
30 Lacombe, T., Cuturi, M. and Oudot, S. 2018. Large scale computation of means and clusters for
31 persistence diagrams using optimal transport. *Advances in Neural Information Processing Systems* 31.
32
33 Larsson, M.N., Morell Miranda, P., Pan, L., Başak Vural, K., Kaptan, D., Rodrigues Soares, A.E., Kivikero,
34 H., Kantanen, J., Somel, M., Özer, F., Johansson A. M., Stora J., Günther T. 2024. Ancient sheep genomes
35 reveal four Millennia of North European short-tailed sheep in the Baltic Sea region. *Genome Biology*
36 *and Evolution* evae114. DOI: 10.1093/gbe/evae114
37
38 Lavocat, R. 1966. *Atlas de Préhistoire. Tome III, Faunes Et Flores Préhistoriques de L'Europe Occidentale*.
39 Boubée et Cie.
40
41 Le Meillour, L., Zazzo, A., Zirah, S., Tombret, O., Barriol, V., Arthur, K.W., Arthur, J.W., Cauliez, J., Chaix,
42 L., Curtis, M.C., Gifford-Gonzalez, D., Gunn, I., Guthertz, X., Hildebrand, E., Khalidi, L., Millet, M.,
43 Mitchell, P., Studer, J., Vila, E., Welker, F., Pleurdeau, D. and Lesur, J. 2023. The name of the game:
44 palaeoproteomics and radiocarbon dates further refine the presence and dispersal of caprines in
45 eastern and southern Africa. *Royal Society Open Science* 10(11): 231002. DOI:
46 <https://doi.org/10.1098/rsos.231002>.
47

1 Le Meillour, L., Zirah, S., Zazzo, A., Cersey, S., Déroit, F., Imalwa, E., Lebon, M., Nankela, A., Tombret,
2 O., Pleurdeau, D. and Lesur, J. 2020. Palaeoproteomics gives new insight into early southern African
3 pastoralism. *Scientific Reports* 10(1): 14427. DOI: <https://doi.org/10.1038/s41598-020-71374-3>.
4
5 Lloveras, L., Rissech, C., Davis, S. and Parés-Casanova, P.M. 2022. Morphological Differences between
6 Sheep and Goat Calcanea Using Two-Dimensional Geometric Morphometrics. *Animals* 12(21): 2945.
7
8 Lundberg, S.M. and Lee, S-I. 2017. A unified approach to interpreting model predictions. *Advances in*
9 *neural information processing systems* 30.
10
11 Lv, F.-H., Cao, Y.-H., Liu, G.-J., Luo, L.-Y., Lu, R., Liu, M.-J., Li, W.-R., Zhou, P., Wang, X.-H., Shen, M., and
12 others. 2022. Whole-genome resequencing of worldwide wild and domestic sheep elucidates genetic
13 diversity, introgression, and agronomically important loci. *Molecular biology and evolution* 39(2):
14 msab353. DOI: 10.1093/molbev/msab353
15
16 Miele, V., Dussert, G., Cucchi, T. and Renaud, S. 2020. Deep learning for species identification of
17 modern and fossil rodent molars. *BioRxiv* 2020–08.
18
19 Moclán., A, Domínguez-García, Á.C., Stoetzel, E., Cucchi, T., Sevilla, P. and Laplana, C. 2023. Machine
20 Learning interspecific identification of mouse first lower molars (genus *Mus* Linnaeus, 1758) and
21 application to fossil remains from the Estrecho Cave (Spain). *Quaternary Science Reviews* 299: 107877.
22
23 Moclán, A., Domínguez-Rodrigo, M. and Yravedra, J. 2019. Classifying agency in bone breakage: an
24 experimental analysis of fracture planes to differentiate between hominin and carnivore dynamic and
25 static loading using machine learning (ML) algorithms. *Archaeological and Anthropological Sciences* 11:
26 4663–4680.
27
28 Mohaseb, A.F., Cornette, R., Zimmermann, M.I., Davoudi, H., Berthon, R., Guintard, C., Cucchi, T.,
29 Hanot, P., Mohandesan, E., Eisenmann, V., Peters J. Mashkour M. 2023. Predictive use of modern
30 reference osteological collections for disentangling the shape of Eurasian equid cheek teeth and
31 metapodials in archaeological material. *International Journal of Osteoarchaeology* 33(5): 938–954.
32 DOI: 10.1002/oa.3255
33
34 Munro, N.D., Bar-Oz, G. and Hill, A.C. 2011. An exploration of character traits and linear measurements
35 for sexing mountain gazelle (*Gazella gazella*) skeletons. *Journal of Archaeological Science* 38(6): 1253–
36 1265.
37
38 Myronenko, A. and Song, X. 2010. Point set registration: Coherent point drift. *IEEE transactions on*
39 *pattern analysis and machine intelligence* 32(12): 2262–2275.
40
41 Nicolau, M., Levine, A.J. and Carlsson, G. 2011. Topology based data analysis identifies a subgroup of
42 breast cancers with a unique mutational profile and excellent survival. *Proceedings of the National*
43 *Academy of Sciences* 108(17): 7265–7270.
44
45 Otter, N., Porter, M.A., Tillmann, U., Grindrod, P. and Harrington, H.A. 2017. A roadmap for the
46 computation of persistent homology. *EPJ Data Science* 6: 1–38.
47

1 Peters, J. 1989. Osteomorphological features of the appendicular skeleton of gazelles, genus *Gazella*
2 Blainville 1816, bohor reedbuck, *redunca redunca* (Pallas, 1767) and bushbuck, *Tragelaphus scriptus*
3 (Pallas, 1766). *Anatomia, Histologia, Embryologia* 18(2): 97–113.
4
5 Pilaar Birch, S.E., Scheu, A., Buckley, M. and Çakırlar, C. 2019. Combined osteomorphological, isotopic,
6 aDNA, and ZooMS analyses of sheep and goat remains from Neolithic Ulucak, Turkey. *Archaeological*
7 *and Anthropological Sciences* 11(5): 1669–1681. DOI: <https://doi.org/10.1007/s12520-018-0624-8>.
8
9 Plummer, T.W., Bishop, L.C. and Hertel, F. 2008. Habitat preference of extant African bovids based on
10 astragalus morphology: operationalizing ecomorphology for palaeoenvironmental reconstruction.
11 *Journal of Archaeological Science* 35(11): 3016–3027. DOI: <https://doi.org/10.1016/j.jas.2008.06.015>.
12
13 Pöllath, N., Alibert, P., Schafberg, R. and Peters, J. 2018. Striking new paths—Distinguishing ancient *Ovis*
14 *orientalis* from its modern domestic descendant (Karakul breed) applying Geometric and traditional
15 Morphometric approaches to the astragalus. In: *Archaeozoology of the Near East XII. Proceedings of*
16 *the 12th International Symposium of the ICAZ Archaeozoology of Southwest Asia and Adjacent Areas*
17 *Working Group, Groningen Institute of Archaeology, June 14-15 2015, University of Groningen, the*
18 *Netherlands*. 2018. pp. 207–225.
19
20 Pöllath, N., Schafberg, R. and Peters, J. 2019. Astragalar morphology: Approaching the cultural
21 trajectories of wild and domestic sheep applying Geometric Morphometrics. *Journal of Archaeological*
22 *Science: Reports* 23: 810–821. DOI: <https://doi.org/10.1016/j.jasrep.2018.12.004>.
23
24 Popkin, P.R.W., Baker, P., Worley, F., Payne, S. and Hammon, A. 2012. The Sheep Project (1):
25 determining skeletal growth, timing of epiphyseal fusion and morphometric variation in unimproved
26 Shetland sheep of known age, sex, castration status and nutrition. *Journal of Archaeological Science*
27 39(6): 1775–1792. DOI: <https://doi.org/10.1016/j.jas.2012.01.018>.
28
29 Prendergast, M.E., Janzen, A., Buckley, M. and Grillo, K.M. 2019. Sorting the sheep from the goats in
30 the Pastoral Neolithic: morphological and biomolecular approaches at Luxmanda, Tanzania.
31 *Archaeological and Anthropological Sciences* 11(6): 3047–3062. DOI: [https://doi.org/10.1007/s12520-](https://doi.org/10.1007/s12520-018-0737-0)
32 [018-0737-0](https://doi.org/10.1007/s12520-018-0737-0).
33
34 Project, TG. 2023. *GUDHI User and Reference Manual*. 3.9.0. GUDHI Editorial Board.
35
36 Prummel, W. and Frisch, H-J. 1986. A guide for the distinction of species, sex and body side in bones of
37 sheep and goat. *Journal of Archaeological Science* 13(6): 567–577. DOI: [https://doi.org/10.1016/0305-](https://doi.org/10.1016/0305-4403(86)90041-5)
38 [4403\(86\)90041-5](https://doi.org/10.1016/0305-4403(86)90041-5).
39
40 Qi, C.R., Yi, L., Su, H. and Guibas, L.J. 2017. Pointnet++: Deep hierarchical feature learning on point sets
41 in a metric space. *Advances in neural information processing systems* 30.
42
43 Rakotomamonjy, A., Bach, F., Canu, S. and Grandvalet, Y. 2008. SimpleMKL. *Journal of Machine*
44 *Learning Research* 9: 2491–2521.
45
46 Ribeiro, M.T., Singh, S. and Guestrin, C. 2016. ‘Why should i trust you?’ Explaining the predictions of
47 any classifier. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge*
48 *discovery and data mining*. 2016. pp. 1135–1144.

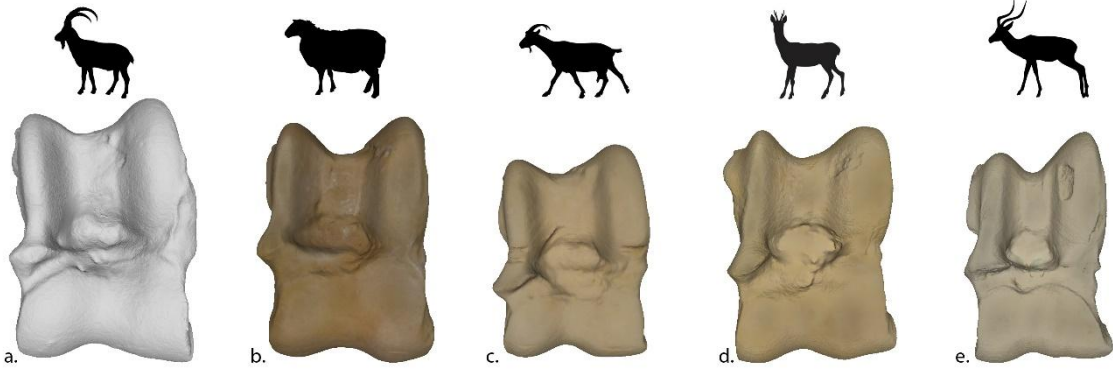
1
2 Rouvreau, V. 2023. Alpha complex. In: *GUDHI User and Reference Manual*. 3.9.0. GUDHI Editorial
3 Board. p.
4
5 Rudin, C. 2019. Stop explaining black box machine learning models for high stakes decisions and use
6 interpretable models instead. *Nature machine intelligence* 1(5): 206–215.
7
8 Salvagno, L. and Albarella, U. 2017. A morphometric system to distinguish sheep and goat postcranial
9 bones Muhlbachler, MC (ed.). *PLOS ONE* 12(6): e0178543. DOI:
10 <https://doi.org/10.1371/journal.pone.0178543>.
11
12 Sipilä, I.M.V., Steele, J., Dickens, L. and Martin, L. 2023. Bones of contention: a double-blind study of
13 experts' ability to classify sheep and goat astragali from images. *Archaeological and Anthropological*
14 *Sciences* 15(12): 187. DOI: <https://doi.org/10.1007/s12520-023-01865-7>.
15
16 Sobolev, A. and Vetrov, D.P. 2019. Importance weighted hierarchical variational inference. *Advances*
17 *in Neural Information Processing Systems* 32
18
19 Tsahar, E., Izhaki, I., Lev-Yadun, S. and Bar-Oz, G. 2009. Distribution and Extinction of Ungulates during
20 the Holocene of the Southern Levant Hansen, DM (ed.). *PLoS ONE* 4(4): e5316. DOI:
21 <https://doi.org/10.1371/journal.pone.0005316>.
22
23 Uerpmann, H-P. 1987. *The ancient distribution of ungulate mammals in the Middle East*. Beihefte zum
24 Tübinger Atlas des Vorderen Orients Reihe A. Naturwissenschaften. Reichert, Wiesbaden.
25
26 Vila, E., Abrahamsi, P., Albesso, M., Amare, A., Bader, C., Berthon, R., Bouzid, S., Bradley, D., Breniquet,
27 C., Chahoud, J., Cucchi, T., Davoudi, H., De Cupere, B., Escarguel, G., Estrada, O., Gourichon, L., Helmer,
28 D., Huangfu, W., Lesur, J., Mashkour, M., Michel, C., Mohaseb, A., Orlando, L., Pompanon, F., Studer,
29 J. and Vuillien, M. 2021. EVOSHEEP: the makeup of sheep breeds in the ancient Near East. *Antiquity*
30 95(379): e2. DOI: <https://doi.org/10.15184/aqy.2020.247>.
31
32 Vuillien, M. 2020. *Systèmes d'élevage et pastoralisme en Provence et dans les Alpes méridionales*
33 *durant la Protohistoire : Nouvelles perspectives en archéozoologie*. Thèse de doctorat de Préhistoire.
34 CEPAM, UMR 7264, Université Côte d'Azur, 687 p.
35
36 Vuillien, M. 2024. *Archéozoologie et Machine Learning : vers une collaboration d'avenir ?* 2024,
37 Bioarcheologies, Le Blog, Available at <https://doi.org/10.58079/VZQ7>.
38
39 Wadsworth, C., Procopio, N., Anderung, C., Carretero, J.-M., Iriarte, E., Valdiosera, C., Elburg, R.,
40 Penkman, K. and Buckley, M. 2017. Comparing ancient DNA survival and proteome content in 69
41 archaeological cattle tooth and bone samples from multiple European sites. *Journal of Proteomics* 158:
42 1–8. DOI: <https://doi.org/10.1016/j.jprot.2017.01.004>.
43
44 Wyatt-Spratt, S. 2022. After the revolution: a review of 3D modelling as a tool for stone artefact
45 analysis. *Journal of Computer Applications in Archaeology* 5(1).
46

1 Zeder, M.A. and Lapham, H.A. 2010. Assessing the reliability of criteria used to identify postcranial
2 bones in sheep, *Ovis*, and goats, *Capra*. *Journal of Archaeological Science* 37(11): 2887–2905. DOI:
3 <https://doi.org/10.1016/j.jas.2010.06.032>.

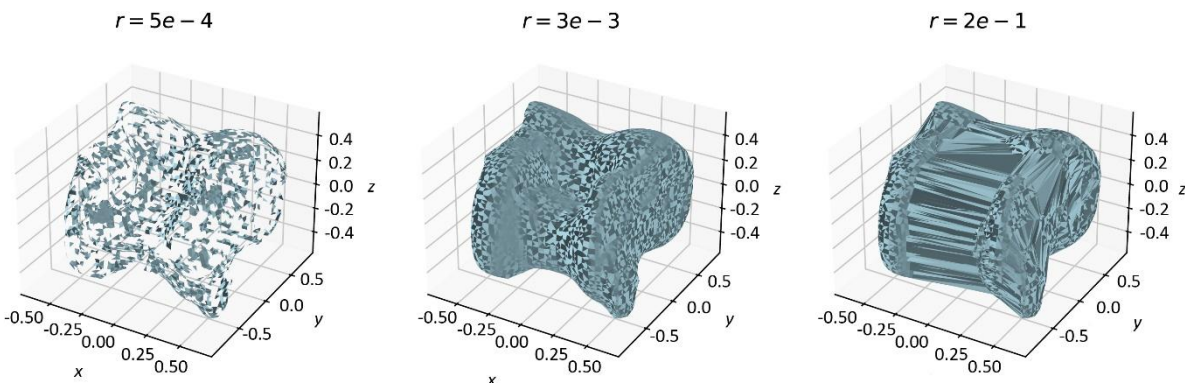
4
5 Zeder, M.A. and Pilaar, S.E. 2009. Assessing the reliability of criteria used to identify mandibles and
6 mandibular teeth in sheep, *Ovis*, and goats, *Capra*. *Journal of Archaeological Science* 37(2): 225–242.
7 DOI: <https://doi.org/10.1016/j.jas.2009.10.002>.

8
9 Zhao, T., Alliez, P., Boubekour, T., Busé, L. and Thiery J.-M. 2021. Progressive discrete domains for
10 implicit surface reconstruction. In: *Computer Graphics Forum*. Wiley Online Library. pp. 143–156. DOI:
11 [10.1111/cgf.14363](https://doi.org/10.1111/cgf.14363)

12
13 Zomorodian, A. and Carlsson, G. 2004. Computing persistent homology. In: *Proceedings of the*
14 *twentieth annual symposium on Computational geometry*. 2004. pp. 347–356.



18
19 **Figure 1:** 3D astragalus presented in dorsal view of (a) Alpine ibex (*Capra ibex*), (b) sheep (*Ovis aries*),
20 (c) goat (*Capra hircus*), (d) roe deer (*Capreolus capreolus*) and (e) gazelle (*Gazella* sp.). Alpine ibex
21 astragalus is untextured.



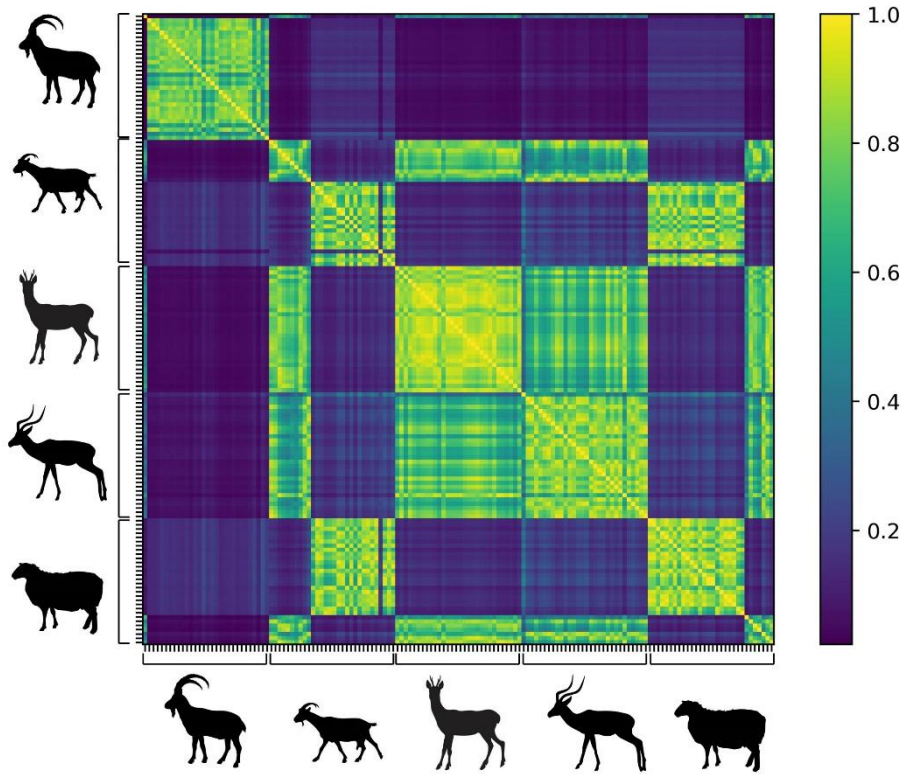
25
26 **Figure 2:** Evolution of the simplicial complex of a 3D astragalus (reference name Obs_1997_187) during
27 the Alpha filtration process, highlighting the impact of varying radius thresholds r . (Left) At this stage,

1 the bone's structure is partially reconstructed, with multiple connected components and some
2 topological cycles, 1-dimensional features, visible. (Center) The Alpha complex has merged into a single
3 connected component, capturing the overall structure of the bone. The cycles from the previous step
4 have been filled (died), and the complex closely approximates the 3D shape of the bone. (Right) At this
5 advanced stage, most topological cycles have disappeared, and the complex over-reconstructs the
6 bone's shape. The filtration process concludes once all topological voids, 2-dimensional features, are
7 filled.

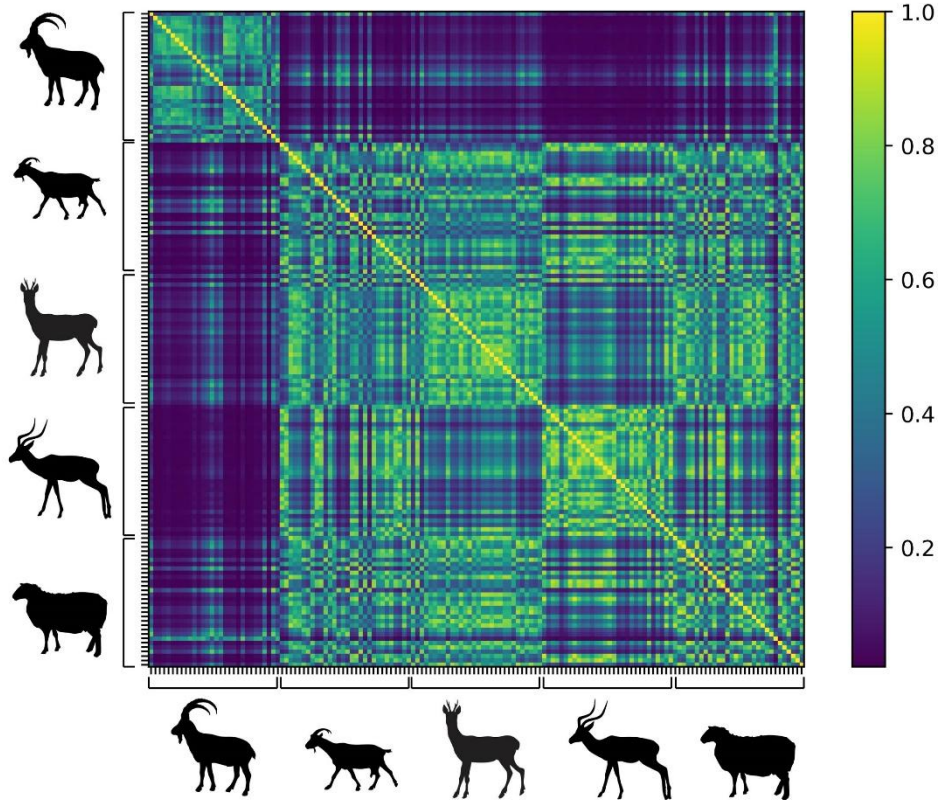
8
9
10
11
12
13
14
15
16

Pre-print document

Wasserstein kernel matrix - dimension 1 without normalisation

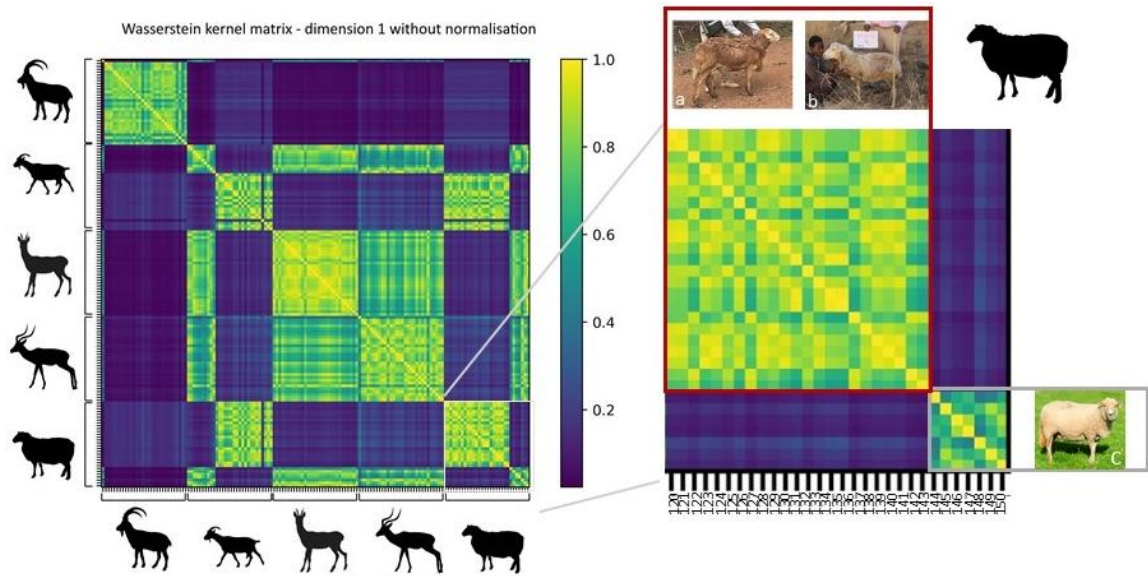


Wasserstein kernel matrix - dimension 2 without normalisation



3 **Figure 3:** Wasserstein kernel matrices without bone's normalisation. Up: topological dimension 1;
4 Down: topological dimension 2. For both the matrices the color code, indicated by the colorbar on the
5 right of the matrix, represents pairwise similarity within the range $[0, 1]$. Yellow cells (similarity equals

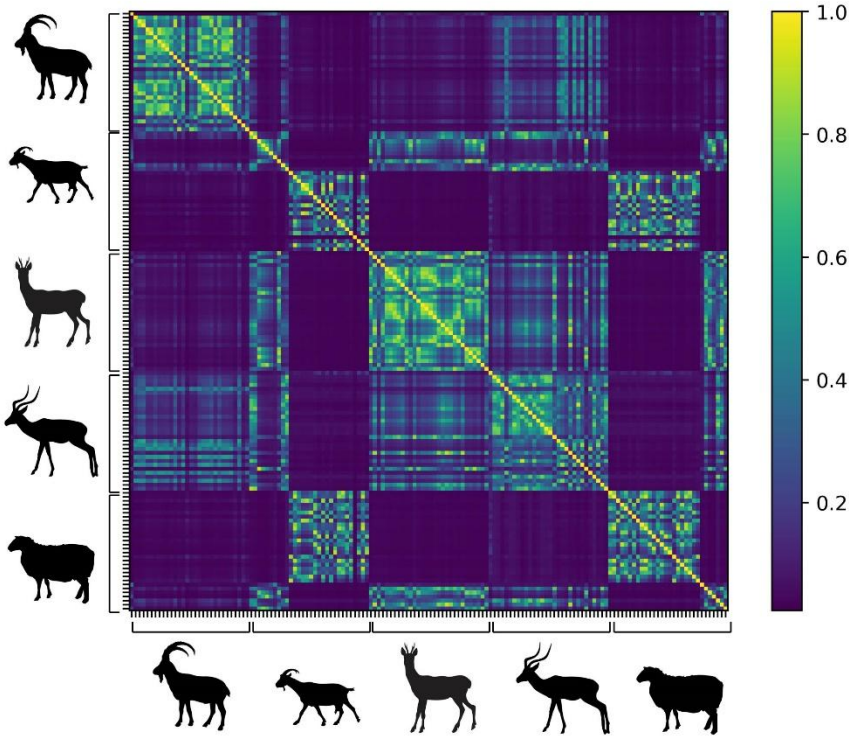
1 to 1), such as those along the diagonal, signify that the x and y bones are identical. As the color shifts
2 towards blue, the bones exhibit increasing dissimilarity (similarity approaching 0).



3
4 **Figure 4:** Illustration of topological dissimilarities observed in sheep in the Wasserstein kernel matrices
5 without bone's normalisation (dimension 1). Picture of sheep breed "Bonga" (a) and "Menz" (b) from
6 Ethiopia © A. Amane / E. Vila. c) Picture of sheep breed "Landes de Bretagne" from France ©H. Ronné
7 <https://www.ecomusee-rennes-metropole.fr/le-mouton-des-landes/>

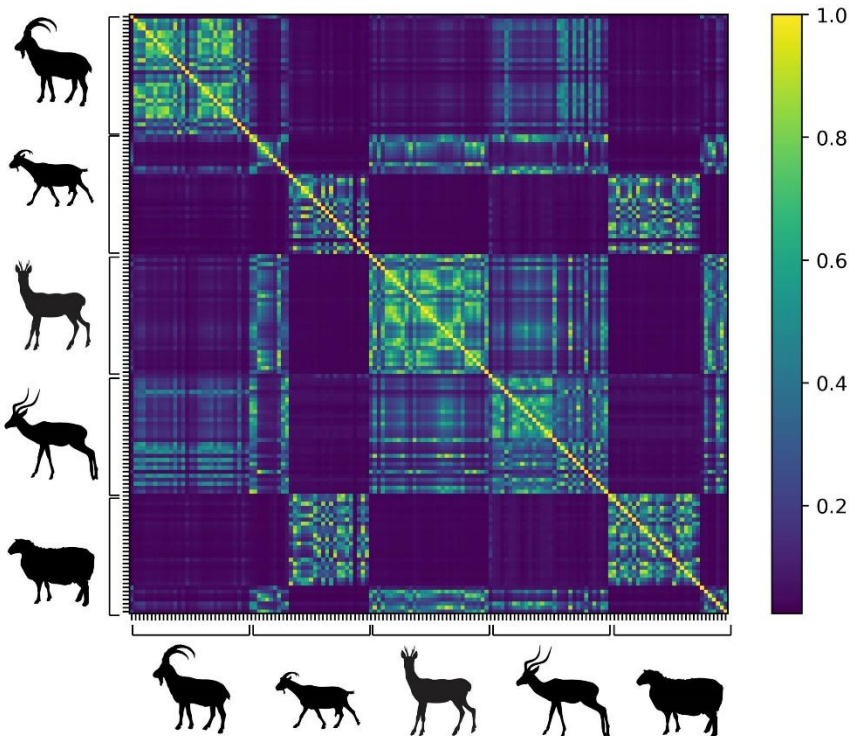
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

Wasserstein kernel matrix - dimension 1 with normalisation



1

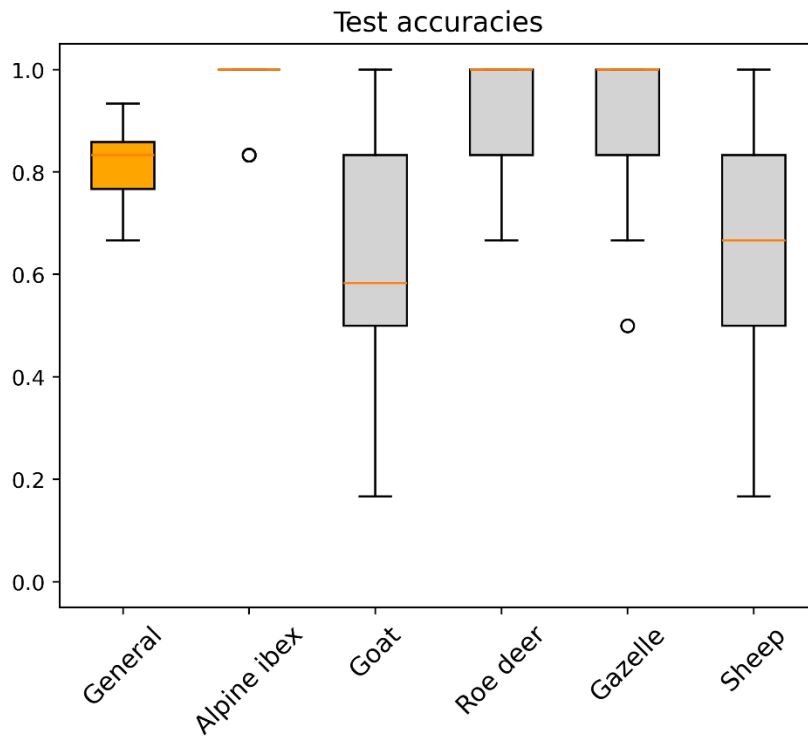
Wasserstein kernel matrix - dimension 1 with normalisation



2

3 **Figure 5:** Wasserstein kernel matrices with bone's normalisation. Up: topological dimension 1; Down:
4 topological dimension 2. For both the matrices the color code, indicated by the colorbar on the right
5 of the matrix, represents pairwise similarity within the range $[0, 1]$. Yellow cells (similarity equals to 1),
6 such as those along the diagonal, signify that the x and y bones are identical. As the color shifts towards
7 blue, the bones exhibit increasing dissimilarity (similarity approaching 0).

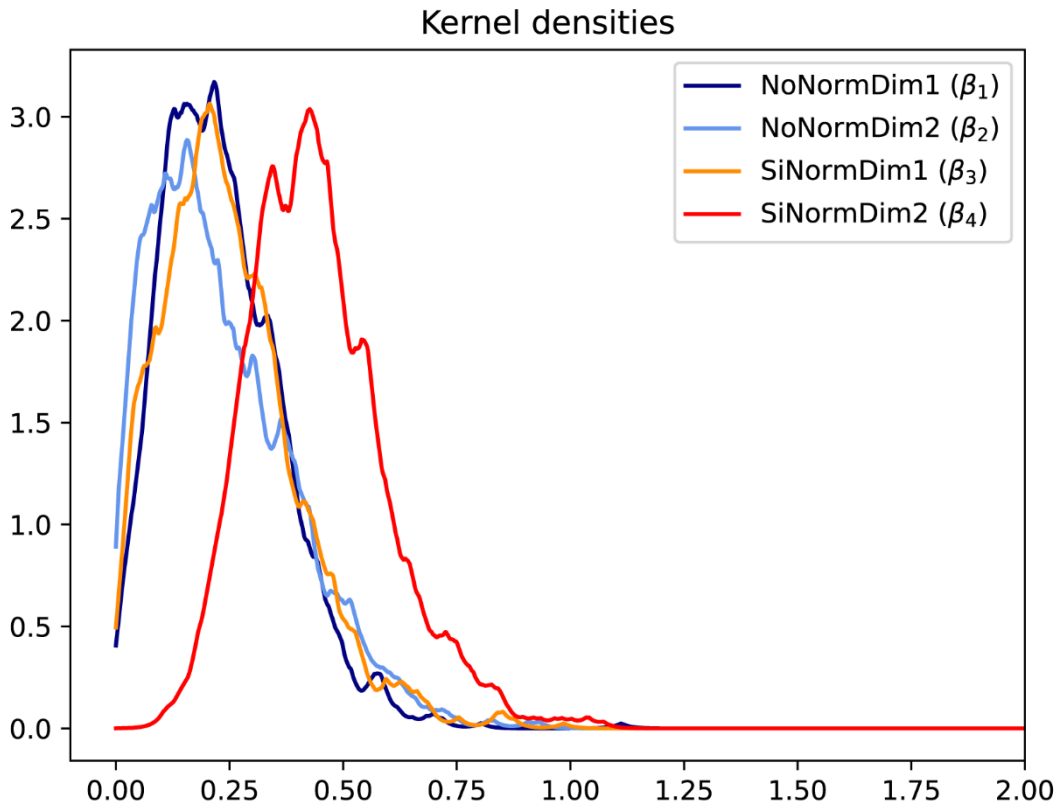
8



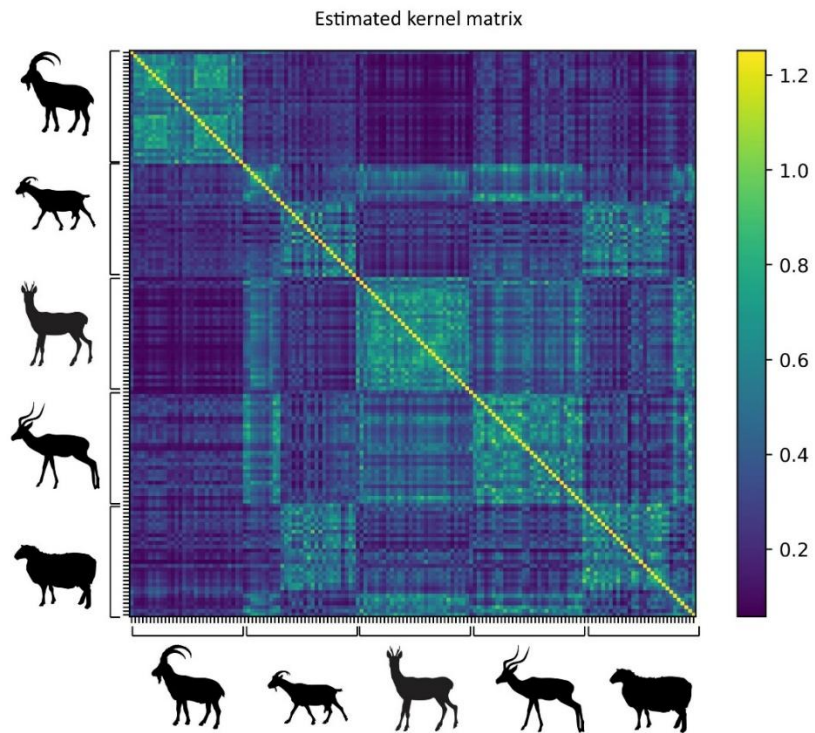
1

2 **Figure 6a** : Boxplot average test classification accuracies. The first column (leftmost) represents the
 3 average performance over the five classes, while the others show the average accuracy for each specie.
 4 The central line (orange) within each box represents the median accuracy, while the lower and upper
 5 edges of the box correspond to the first (Q1) and third quartiles (Q3), respectively, indicating the
 6 interquartile range (IQR). The whiskers extend to the minimum and maximum values within (1.5*IQR)
 7 from Q1 and Q3. Beyond this range are considered outliers and are shown as individual markers.

8



1
 2 **Figure 6b** : Kernel densities estimates. Different colors correspond to the weights of the four
 3 Wasserstein kernel matrices obtained via TDA. X-axis corresponds to weights values, while the y-axis
 4 indicates the estimated density, reflecting the relative frequency of occurrence. Peaks in the density
 5 curves show the most probable values of the weights, while the spread of each distribution provides
 6 insight into the variability and uncertainty of the estimates.
 7



1
 2 **Figure 7:** Learned kernel matrix for a single data split. It represents the optimal linear combination,
 3 where each input kernel matrix is weighted by its corresponding estimated weight value. The color
 4 code follows the same interpretation as in Figures 3 and 5.
 5