

Artificial Intelligence in Radiolarian Fossil Identification; Taxonomic, Biostratigraphic and Evolutionary implications

Veronica Carlsson

► To cite this version:

Veronica Carlsson. Artificial Intelligence in Radiolarian Fossil Identification; Taxonomic, Biostratigraphic and Evolutionary implications. Paleontology. Université de Lille, 2023. English. NNT: . tel-04471671v2

HAL Id: tel-04471671 https://theses.hal.science/tel-04471671v2

Submitted on 7 Mar 2024

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 4.0 International License







Thèse de doctorat

Université de Lille Sciences de la Terre et l'Univers

Veronica Carlsson

CNRS UMR 8198 – Evo-Eco-Paléo & CNRS UMR 9189 – CRIStAL

Artificial Intelligence in Radiolarian Fossil Identification

Taxonomic, Biostratigraphic and Evolutionary implications

Détermination des Radiolaires Fossiles à l'aide de l'Intelligence Artificielle

Implications Taxonomiques, Biostratigraphiques et Évolutives

Pr. Taniel Danelian Université de Lille, France Directeur de thése Pr. Pierre Boulet Université de Lille, France Co-directeur de thése Pr. Cordey Fabrice Université de Lyon 1 Rapporteur Dr. Thibault de Garidel-Thoron CNRS, Université d'Aix-Marseille Rapporteur Pr. François Danneville Université de Lille, France Président du jury Dr. Allison Hsiang Université de Stockholm, Suède Examinatrice Pr. Rie Hori Université de Ehime, Japon Examinatrice Pr. Catherine Crônier Université de Lille, France Examinatrice Dr. Philippe Devienne CNRS, Université de Lille Invité

Membres du jury de la soutenance de thèse de Veronica Carlsson le 8 décembre 2023, à Lille :



Acknowledgments	. 7
CHAPTER 1: Introduction	. 8
1.1 The essence of this research	. 8
1.1.1 The middle Eocene	. 9
1.1.2 What are radiolaria?	. 9
1.1.3 From living organisms to several million-year-old preserved fossils	11
1.1.4 Taxonomic importance of radiolaria	12
1.1.5 Radiolaria used in biostratigraphy	12
1.1.6 Geologic setting of ODP Leg 207	15
1.1.7 The importance of Radiolaria for paleoenvironments and paleoceanography	17
1.1.8 Machine learning techniques and Artificial Intelligence	18
1.1.9 Earlier research using AI in image recognition in micropaleontology	20
1.2 Thesis objectives and structure	21
PART 1: ARTIFICIAL INTELLIGENCE APPLIED TO INDIVIDUAL SPECIES	24
CHAPTER 2: Artificial intelligence applied to the classification of eight middle Eocene species	
of the genus <i>Podocyrtis</i> (Polycystine radiolaria)	24
Author contributions	25
Abstract	25
2.1 Introduction	26
2.2 Materials	30
2.3 Methods	37
2.3.1 Manual picking of individual <i>Podocyrtis</i> specimens	37
2.3.2 Image acquisition and processing	38
2.3.3 Datasets	40
2.3.4 MobileNet convolutional neural networks	41
2.4 Results	43
2.4.1 CNN accuracies	43
2.4.2 Confusion matrices	44
2.4.3 Testing the predictive models	49
2.5 Discussion	50
2.5.1 MobileNet performance and accuracy	50
2.5.2 Species choice and their image properties	55
2.6 Conclusions	56
Data availability	58

Supplement	58
Acknowledgements	58
CHAPTER 3: Morphometrics and machine learning discrimination of the middle Eocene	e
radiolarian species <i>Podocyrtis chalara</i> , <i>Podocyrtis goetheana</i> and their morphological intermediates	60
Author contributions	61
Abstract	62
3.1 Introduction	63
3.2 Analyzed morphological groups	66
3.3 Materials and methods	
3.3.1 Sediment samples	
3.3.2 Slide preparation	
3.3.3 Microscopy and image processing	
3.3.4 Morphometric analyses	
3.3.5 Artificial Neural Networks	
3.4 Results	81
3.4.1 Morphometrics and linear discriminant analysis	81
3.4.2 Artificial neural networks	
3.4.2.1 Classification using the four morphogroups	83
3.4.2.2 Classification using the three morphogroups supported by LDA	85
3.5 Discussion	88
3.6 Conclusion	
Data availability	
Acknowledgments	
PART 2: ARTIFICIAL INTELLIGENCE APPLIED TO ASSEMBLAGES	
CHAPTER 4: Subchapter 1 - Convolutional neural network application on a new middle	Eocene
radiolarian dataset	
Author contributions	
Abstract	
4.1 Introduction	
4.2 Materials and methods	
4.2.1 Core setting and sample preparation	
4.2.2 Image collection and processing	100
4.2.3 Taxa selection and dataset	102
4.2.4 CNN training	107

4.2.5 Performance metrics	107
4.2.6. Test set to validate the CNN	
4.2.7 Application on 39 species	
4.3 Results	109
4.3.1 Training of the initial dataset	109
4.3.2 Performance validation from the test set	
4.3.3 Application on new samples	
4.4 Discussion	
4.4.1 Classification	
4.4.2 CNN training and new test set score	117
4.4.3 Feedback on individual species	117
4.5 Conclusions	
Data availability	
Acknowledgments	
Appendix A. Supplementary data	
CHAPTER 4: Subchapter 2 - Initial middle Eocene radiolarians dataset from the tro Atlantic (ODP Leg 207) partly classified by K-means clustering	opical 128
4.1 (2) Introduction	
4.2 (2) Results	
4.3 (2) Discussion	
Appendix A: Initial list of all the classes for the radiolarian image database from the De at the Tropical Atlantic Ocean.	emerara Rise 132
CHAPTER 5: Using convolutional neural network for improving the biostratigraphy Eocene radiolaria from ODP Leg 207. Site 1260	y on middle 136
Abstract	
5.1 Introduction	
5.2 Materials and Methods	
5.2.1 Materials	
5.2.1 Materials	
 5.2.1 Materials 5.2.2 Methods 5.3 Results 	
 5.2.1 Materials 5.2.2 Methods 5.3 Results	
 5.2.1 Materials 5.2.2 Methods 5.3 Results	
 5.2.1 Materials	
 5.2.1 Materials	

5.4.1 General discussion	150
5.4.2 Individual species revision and their credibility in future biostratigraphic work	150
5.4.3 Time and Effort	155
5.5 Conclusion	155
Acknowledgments	156
CHAPTER 6: General Conclusions and Perspectives	157
6.1 General conclusions	157
6.2 Perspectives	161
References	165
Appendices	184
Abstract	185
Résumé	278

Acknowledgments

I would first like to thank the Marie Curie Skłodowska grant and the PEARL program for giving support, not just in form of salary and scholarship, but also providing help when it comes to anything from problem with the accommodation and other administrative things, also for providing the absolute best programs and encouraging us about options and how to proceed in one's career both inside and outside of academia. Also, a huge thanks for the other international students of the PEARL program, that have been a huge help from handling any kind of problems or situations normally faced for expats, especially when not being fluent in french.

Thanks to my supervisors Taniel Danelian, Pierre Boulet and Philippe Devienne for guiding me throughout the project, and thanks to both the Evo-Eco-Paléo lab, the CRIStAL lab, especially the neuromorphic group at IRCICA, for their regular meetings. Thanks also to my other co-authors present at the same or other labs or institutions around the world. A huge thanks also to the master students, Aurelien Laforge and Jose Francisco Pinto Cabrera, for their excellent contribution to this work.

A massive thanks to Johan Renaudie that was able to host me at Museum für Naturkunde in Berlin and to Allison Hsiang for hosting me at Stockholm University, for a certain period of time during this PhD.

I would like to give special thanks to my family and my husband, who has always been there for me from a distance, giving me support, energy and motivation to keep on going and to my family's close friends in Lille that has helped and supported me no matter the circumstances.

And at last, thank you Ernst for everything, you are forever missed and my biggest inspiration.

CHAPTER 1: Introduction

1.1 The essence of this research

Paleontological data, i.e., fossil remains of dead organisms, can give important understandings on changes in paleodiversity throughout geologic time. The study of past environments and climates and the biotic response to these events are studied by paleontological data, especially with planktonic microfossils. Plankton are small single celled microorganisms living in aquatic environments. They are easy to study due to their small size, high abundance and rapid evolution and can be studied by just extracting them from very small samples. This also enables them to be easily used for biostratigraphy i.e., dating of sediments based on species occurrences. For being able to make any type of interpretations on diversity changes, biostratigraphy, or other abundance data studies, a huge amount of paleontological data must be collected and analysed which is rather tedious, exhausting and time consuming, but yet a very important work.

The aim of this study is both simplifying the work of fellow taxonomic experts in radiolarian taxonomy by tackling unagreed taxonomic challenges and keep a consistency. It also aims to simplify or automatize the work, so even non-taxonomic experts can use AI for automatic image recognition. This type of research could enhance directly or indirectly the scientific field by automatizing or simplifying the research about radiolarians, their important role with respect to dating sediments used for biostratigraphy in areas where the better studied calcareous nannofossils or foraminifera are not enough. Also, it would be useful for paleoclimatic or paleoenvironmental studies, as well as to study the general response of plankton to long- or short-term climatic events. In industry, this research is very useful for enhancing or partly replacing the work made by biostratigraphers in oil and gas exploration, who use microfossils

to date sediments and to recognize different source rocks, which is a vital part of finding oil and gas. Here we focus on radiolarians coming from the middle Eocene time.

1.1.1 The middle Eocene

The Eocene is a time period stretching as far back as 56-33.9 Million years ago (Ma). At the beginning of this period, both Australia and Antarctica were still connected. Also, more major tectonic events occurred during this interval, such as the disappearance of the Tethys Ocean, leading to the formation of the Alpes following the collision of India with Asia. The Eocene is also characterized by significant changes in long-term climate change, with temperatures reaching the highest levels of the Cenozoic at the beginning of the period, which is known as the Paleocene-Eocene Thermal Maximum (PETM). After this peak, the overall long-term temperature followed a global cooling trend, which continued throughout the entire Cenozoic until present day. The middle Eocene, which is the interval of radiolarians studied in this work, is characterized by a long-term cooling, and a single warming event, the Middle Eocene Climatic Optimum (MECO).

1.1.2 What are radiolaria?

Radiolaria are a group of single-celled eukaryotes, Protozoa, living in aquatic environments; they are found in all oceans and at all ranges of depths, from the surface to the deep sea. They are very diverse today and throughout geologic time and are very important components of the marine ecosystem, serving as prey for larger animals such as filter-feeding invertebrates; they feed on a variety of smaller organisms, like phytoplankton (Anderson, 1983; Boltovskoy et al., 2017; and Matsuoka, 2007). They are important players in both the carbon (Lampitt et al., 2009)

and silica cycles (Takahashi, 1987). Radiolaria are closely related to foraminifera (Cavalier-Smith et al., 2018) and they both belong to the eukaryotic supergroup Rhizaria.

The radiolaria that are of most interest to paleontology are Polycystines, both because they have been abundant throughout geological time, but also since their test are made of opaline silica which is rather well preserved in the fossil record. They have generally a spherical or cylindrical shape and complex skeletons, which they use for support and protection. Two of the Polycystines orders are Spumellaria and Nassellaria (Fig. 1.1).

Spumellarians are spherical or ovoid with a radial symmetry and bear a complex, often spiny, external skeleton and one or several inner shells. Spumellarians have existed since the Ordovician, approximately 475 Ma ago. Living Spumellaria have been found to live in symbiosis with several different photosynthetic algae, for instance, dinoflagellates, cyanobacteria, prasinophytes, or haptophytes. This symbiotic association is likely the reason of the higher concentrations and greater diversity of Spumellaria in the tropics (Sandin, 2019).

Nassellarians, unlike Spumellarians, have existed since the Carboniferous, approximately 360-300 Ma. The skeleton of Nassellarians is typically cone- to cylindrical-shaped and it has a bilateral symmetry. It is assumed that species diversity in Nassellarians has been increasing since the beginning of the Cenozoic, approximately 66 Ma (Anderson, 1983).



Diagrammatic idealised nassellarian radiolaria showing basic morphological features

Figure 1.1. Typical sketches of Spumellarian and Nassellarian skeletons including names for the different skeletal parts. Slightly adjusted from Anderson, (1983).

1.1.3 From living organisms to several million-year-old preserved fossils

Once radiolaria die, only their siliceous skeleton remains and sinks to the bottom of the ocean and accumulates slowly into the sediment and will eventually through time undergo diagenesis; thus, the sediment in the radiolarian skeleton remains will undergo pressure and chemical changes and will lithify (become rock). Biological imprints will eventually get fossilized and become excellent sources for reconstructing past environments (De Wever et al., 2002). Polycystine radiolaria are well preserved throughout the fossil record and are unique among all microfossil groups by having a complete fossil record with a detailed evolutionary history (Sanfilippo and Riedel, 1990).

1.1.4 Taxonomic importance of radiolaria

Generally, the classification of extinct organisms is different from biology. In paleontology, classification is based on morphological features. Whereas, in biology, classification on a species level is based on the possibility of producing fertile offspring.

Several hundreds or thousands of extinct or extant radiolaria species have been recognized. Many of these are still not discovered or described by the literature. Generally speaking, the radiolarian taxonomy is still not fully understood due to difficulties in observing the internal parts of the skeleton in radiolarians. Many radiolarian species are hard to place in higher ranks such as genera or families, and many species can be placed into many different families. This is mainly because there are not enough phylogenetic traces for species in each genus (Suzuki and Aita., 2010). Therefore, the classification of radiolaria based on a species level is very important. Haeckel was one of the first people trying to organize and classify radiolarians into higher taxonomic ranks, and there have been several other more or less failed attempts to group radiolarians (De Wever et al., 2001).

Lately, there have been attempts to re-classify certain species into other genera or families by using and connecting molecular studies to extinct species based on extant taxa (Suzuki et al., 2021; Sandin, 2019) but still upon today the Haeckel taxonomy is still dominating within radiolaria research, especially for Cenozoic radiolaria (Lazarus, 2005).

1.1.5 Radiolaria used in biostratigraphy

Many microscopic unicellular organisms with fossilized skeleton parts are thanks to their fast evolution, high abundance, and small size extremely useful to biostratigraphy. Fossils themselves cannot be dated using radiometric dating techniques. There are however other proxies that can be correlated with fossil dating. For example, paleomagnetism is based on the knowledge that the magnetic north and south poles are constantly moving creating changes in polarity and by studying the alignment of the magnetic mineral crystallization in for example magnetite or iron, changes over time can be recorded. This is correlated with well-dated reference data. Tephrochronology is another age correlation technique. The tephra layers consisting of volcanic ash deposits in oceanic sediments can be used to correlate age and compare it with known ages for volcanic eruptions, which have been well-dated.

Originally, the most studied microfossil groups, which are widely used for biostratigraphy, including for oil and gas exploration, are foraminifera, nannofossils, and palynomorphs. Cenozoic biozonataions based on calcareous nannoplankton and planktonic foraminifera, have approximate ranges for 2 million years (Bolli et al., 1985; Armstrong and Brasier, 2005).

It was earlier believed that radiolarian species, initially described by Haeckel, were more longranging and could therefore not be useful within biostratigraphy (Campbell, 1984; Lazarus, 2005). However, this was not the case and radiolarians evolved at a similar rate as other microorganisms and radiolarians are therefore useful in areas lacking or having too poorly preserved carbonate microfossils (Lazarus, 2005).

In 1985, Sanfilippo et al. (in Bolli et al., 1985) conducted an important work in reviewing 29 existing radiolarian biozones correlated with magnetostratigraphy (Armstrong and Braiser, 2005).

Sanfilippo and Nigrini (1998) revised many of the tropical Eocene biozonations used today in biostratigraphy. Many of the *Podocyrtis* species are used as biomarkers for the different radiolarian zonations along with a few other Eocene radiolarian species. The genus *Podocyrtis*

is generally well studied, and its different species occur during the middle Eocene. The first time this genus was described was by Ehrenberg in 1847. Sanfilippo and Riedel conducted a morphometric study of the evolution of the *Podocyrtis* species in 1990, showing the *Podocyrtis* species evolving during the middle Eocene time frame.

The radiolarian biozonations have been calibrated by geomagnetic polarity timescales (Kamikuri et al., 2012; Souza et al., 2017, 2021; Hollis et al., 2020). The important biozonations covering our studied materials are from:

RP16 – *Podocyrtis (Lampterium) goetheana* interval zone (Moore, 1971, emend. Riedel and Sanfilippo, 1978). Its base is recognized by the First Occurrence (FO) of *Podocyrtis (Lampterium) goetheana* and the top is identified by the FO of *Cryptocarpium azyx* (Sanfilippo and Riedel).

RP15 – *Podocyrtis (Lampterium) chalara* lineage zone (Riedel and Sanfilippo, 1970, 1978). The base is distinguished by the Evolutionary Transition (ET) of *Podocyrtis (Lampterium) mitra* Ehrenberg to *Podocyrtis (Lampterium) chalara* Riedel and Sanfilippo.

RP14 – *Podocyrtis (Lampterium) mitra* lineage zone (Riedel and Sanfilippo, 1970, 1978). The base is recognized by the ET of *Podocyrtis (Lampterium) sinuosa* Ehrenberg, to *Podocyrtis (Lampterium) mitra* Ehrenberg.

RP13 – *Podocyrtis (Podocyrtoges) ampla* lineage zone (Riedel and Sanfilippo, 1970, 1978). The base is set by the ET of *Podocyrtis (Podocyrtoges) physis* Sanfilippo and Riedel, to *Podocyrtis (Podocyrtoges) ampla* Ehrenberg. RP12 – *Thyrsocyrtis (Pentalocorys) triacantha* lineage zone (Riedel and Sanfilippo, 1970; emend. Riedel and Sanfilippo, 1978). Which the base is based on the FO of *Eusyringium lagena* (Ehrenberg).

A biostratigraphic research based on radiolarians (Meunier and Danelian, 2022) was recently conducted from the samples used in this work which were correlated with magnetostratigraphy and cyclostratigraphy from site 1260 (Westerhold and Röhl, 2013), resulting in a highly accurate and high resolutional ages for up to 71 different radiolarian bioevents.

1.1.6 Geologic setting of ODP Leg 207

The location of the studied radiolarian samples is from the Demerara Rise, a plateau, off the Atlantic coast of Suriname and French Guyana in South America (Fig. 1.2 and 1.3). Sediment samples were collected there from the Ocean Drilling Program (ODP) expedition Leg 207 in 2004 with the aim to recover extended shallowly submerged Cretaceous and Paleogene sediments that could be used for palaeoceanographic studies of the tropical Atlantic Ocean. It was previously known that during this period there were a lot of events happening like periods of ocean anoxia, fast changes in the climate, mass extinction events, and the opening of the Equatorial Atlantic Gateway (Shipboard Scientific Party, 2004). When the sediment samples had been recovered it was discovered that very well-preserved and completed sequences of radiolaria and other siliceous micro remains exist, indicating silica-enriched waters. Hence, the radiolaria diversity and abundance at the Demerara Rise was particularly huge, especially during the middle Eocene. The sedimentation rate during the middle Eocene was also very high, which made it possible for high-resolution studies from this sequence. Due to these facts, this sequence is a great choice for studying Middle Eocene radiolarians.



40 Ma Reconstruction

Figure1.2.PaleogeographicWorldmapfrom40Ma.Constructedbyhttps://www.odsn.de/odsn/services/paleomap/paleomap.html#Form(Hay et al., 1999).

Table 1.1. Samples from ODP Leg 207 from the Demerara rise and their exact coordinates, water depths and sedimentation depths.

Leg/Hole	Coordinates	Water depth (mbsl)	Sediment depth (mbsf)
207-1258A	9°26.00030′N, 54°43.9994′W	3192 mbsl	8-253 mbsf
207-1259A	9°17.99890′N, 54°11.9984′W	2354 mbsl	125-445 mbsf
207-1260A	9°15.94850'N, 54°32.6327'W	2549 mbsl	35-335 mbsf



Figure 1.3. Sample location (Shipboard Scientific Party, 2004)

At this time, the location of the Demerara Rise, located in the northeastern part of the South American continent, was much closer to the African continent than it is today. North and South America were also not connected, leading to potential influences of Pacific Ocean waters, and the Tethys Sea north of Africa.

1.1.7 The importance of Radiolaria for paleoenvironments and paleoceanography

Throughout geologic time, there have been major changes in the climate, with drastic and more long-term changes in temperatures. Understanding the biotic response to these events could help in understanding the ecological and biological effects of the drastic global warming that the Earth is facing today. Understanding the response in microorganisms, which are the primary producers and first-order consumers, is crucial. As of today, the response of radiolaria to global change is not well known, especially their response to diversity or abundance. As mentioned in the taxonomic part, there are a lot of difficulties dealing with taxonomy and we are far from all species being described in the literature, making it difficult to work on reliable diversity or correct abundance data.

This reflects the current situation of our samples, as we are dealing with tropical Atlantic radiolarian assemblages and the middle Eocene radiolarian diversity is huge and based on our estimates it accounts for ca. 500 species, many of which are not described yet.

1.1.8 Machine learning techniques and Artificial Intelligence

There are many different morphometric tools and methods that are frequently used in paleontology to discriminate closely related or similar-looking species or morphotypes based on morphological criteria. One supervised method to find similar data points that can be used to discriminate morphospecies among labeled groups is Linear Discrimination Analysis (LDA). By using a set of different measured characteristics for labeled individuals, LDA can use this information to reduce the dimensions and find the common denominator for each group.

K-Nearest Neighbour (KNN) is another supervised machine learning technique that works by only feeding in images and letting the algorithm find common data points (pixel values) within the different groups.

During the last couple of years, Artificial Intelligence (AI) has become more advanced and developed in a way that can improve these problematic classification tasks. Convolutional Neural Networks (CNNs) are a class within AI, that are popular to use for analysis and classification of images. The function behind CNNs is that it recognizes patterns in images. The CNNs are simply constructed by input layers, which transfer their information into hidden layers, which are the convolutional layers that will process and transform the information into the next layers. Each convolutional layer has different sizes, and numbers of filters. A filter can be seen as a small grid of pixels with different pixel values in each grid corresponding to a specific colour value. This grid will go through the entire image in a sliding (convolving) way and transform the new values to the next layer that will process the image in a different or similar way. Early layers could for example easily detect edges, circles, or corners, and later layers can even recognize more specific objects.

Spiking Neural Networks (SNNs) is a third-generation neural network and is constructed in a way that closely imitates biological brains. In such networks, information is encoded as streams of time events, the spikes, and processes these spike streams in an asynchronous event-based way instead of in a fixed clock cycle. A SNN uses analog neurons, which in comparison with CNN do not fire at each propagation cycle. A neuron fires (produces a spike) when it has reached a certain value, a "threshold". It takes time for neurons to fire and to return to a stable state after being fired. The time after a neuron has spiked until it has returned to its stable state is called the "refractory period", meaning that during this time, the neuron does not respond to any input activity. Learning in SNNs can be supervised (requiring a labeled training dataset) or unsupervised (allowing automatic clustering of the samples of the training dataset). Another strength of SNNs is their ultra-low power consumption when implemented in hardware, opening the way in the future to autonomous automatic classifying microscopes. On the negative side, they are not as well understood as other Artificial Neural Networks (ANNs) and their classification performance lags behind that of ANNs.

1.1.9 Earlier research using AI in image recognition in micropaleontology

There have been several studies trying to apply automatic image recognition using AI or more specifically CNNs. Already in 1996, the first automatic image recognition was used for coccolithophores. The SYstème de Reconnaissance Automatique de Coccolithes (SYRACO), version 1 was presented at a conference, obtaining an average classification accuracy of 49 %. In 1999, a second version of SYRACO was published in a paper by Dollfus and Beaufort, using a simple four-layered deep neural network, which obtained a mean accuracy of about 86%. It wasn't until two years later that the SYRACO was implemented, for the first time, into a real scientific approach for reconstructing Pleistocene oceanic primary productivity (Beaufort et al., 2001). Throughout the years, SYRACO has continued to be developed and so has a lot of other AI studies for automatic image recognition of other microfossil groups.

This thesis has been inspired by some pioneering earlier works that used automatic image classification of radiolarians, such as the studies conducted by Renaudie et al. (2018) and Tetard et al. (2020). Renaudie et al. (2018) applied the MobileNet CNN architecture for the automatic radiolarian classification of 16 closely related species of the Cenozoic genera *Antarctissa* and *Cycladophora*. This work resulted in a classification accuracy of about 73 %, which was increased to ca. 90 % after ignoring unclassified specimens. Tetard et al (2020) described a new workflow for radiolarian image acquisition, including sediment preparation from the removal of non-siliceous particles, to slide settlements on small 1.2 x 1.2 cm large coverslips on decanters built with a 3D printer to help with a random and uniform settling of particles and also a 3D printed coverslip guide, to ease the application of several small coverslips into one single glass slide. This work also trained a ResNet50 architecture CNN on 132 classes of

Miocene to recent radiolarians and obtained an accuracy of about 90 %, and all work was done with the software ParticleTrieur.

Several other important studies are using automatic image recognition for fossil identification. A few of them are on foraminifera (ex. Mitra et al., 2019, Hsiang et al., 2019; Carvalho et al., 2020; De Lima et al., 2020; Marchant et al., 2020), pollen (ex. Goncalves et al., 2016; Bourel et al., 2020) and even some other radiolarian studies (ex. Itaki et al., 2020).

1.2 Thesis objectives and structure

The general objective of this thesis is to apply Artificial Intelligence for automatic image recognition of tropical Atlantic middle Eocene radiolarian taxa.

The main scientific questions that we will try to address in this thesis are as follows:

Q.1 How can neural network learning achieve equal accuracy in the identification of middle Eocene radiolarian species as an Eocene specialist in radiolarian taxonomy?

Q.2 How can machine learning techniques be compared to classical morphometric analyses?

Q.3 How can Spiking Neural Networks (SNNs) be compared with other deep learning methods?

Q.4 Can a trained neural network on a whole middle Eocene tropical radiolarian assemblage be directly applied for example to biostratigraphic or palaeoceanographic studies?

In order to address these questions two distinct approaches were followed; for chapters 2 and 3 (Part 1) the images used were taken on individually picked specimens, while for chapters 4 and

5 (Part 2), the image dataset used was obtained after screening the entire radiolarian assemblage with the use of an automated microscope.

In Chapter 2, we tackle the first scientific question (Q.1) of this thesis; focused on eight closely related species belonging to the genus *Podocyrtis*.

The first step here was to train a labelled convolutional neural network, in which we had to sample a lot of images for each species, for which specimens were picked up individually; we thus formed a well selected database of images, which were later trained on a MobileNet v1 CNN architecture and obtained a high training accuracy for automatic image recognition. At a later stage, we wanted to test the trained CNN, to see if the trained automatic image classification works in a real case study. For this test, we obtained more radiolarian images from other parts of the ocean of the same genus and let the CNN decide the label. This was later checked on, and in many cases, the CNN could identify the correct specimen or at least classify it to a neighboring species along a lineage.

Chapter 3 represents a follow-up study on *Podocyrtis* species, where we attempt to address the presence of intermediate forms between species *P. chalara* and *P. goethean*a, but also the two last scientific questions (Q.2 and Q.3). We tried to use supervised morphometrical analyses, such as Linear Discrimination Analyses (LDA), to investigate if each individual form shares or has their morphospace. We also used the morphospecies along with intermediate forms to train in ANNs. Except for using CNNs, we also introduced the use of SNNs. Both AI and morphometrics showed that the possible intermediate form resembling to *P. chalara* is in fact only recognized as *P. chalara*; on the contrary, the intermediate form closer to *P. goetheana*

by the morphological disparity analysis and by the AI. From the different types of AI analyses, we could conclude that using a SNN is faster and is known to be of less energy consumption than using a CNN.

In both Chapters 4 and 5, we deal with the last scientific question (Q.4). In Chapter 4 we wanted to deal with the whole middle Eocene radiolarian assemblages but focused only on 39 important Nassellarian species. Here, we trained all objects and radiolarians appearing, preferably at the species level; however, this was not always possible and therefore higher taxonomic ranks were also included. We obtained a high training accuracy and applied this trained CNN to new samples, in which we wanted to focus on finding a possible application that could be easily obtainable with this specific trained neural network.

In Chapter 5, we focus only on the automatic image application that is solely based on the trained CNN of Chapter 4. Here we applied biostratigraphy, only focusing on the biostratigraphically important species.

PART 1: ARTIFICIAL INTELLIGENCE APPLIED TO INDIVIDUAL SPECIES

<u>CHAPTER 2</u>: Artificial intelligence applied to the classification of eight middle Eocene species of the genus *Podocyrtis* (Polycystine radiolaria)

Veronica Carlsson^{1,2}*, Taniel Danelian¹, Pierre Boulet², Philippe Devienne², Aurelien Laforge^{1,2} and Johan Renaudie³

¹Univ. Lille, CNRS, UMR 8198, Evo-Eco-Paleo, F-59000 Lille, France.

²Univ. Lille, CNRS, CRIStAL – Centre de Recherche en Informatique Signal et Automatique de Lille, UMR 9189, F-59000 Lille, France.

³Museum für Naturkunde, Leibniz-Institut für Evolutions- und Biodiversitätsforschung, 10115 Berlin, Germany

* Corresponding author. E-mail address: veronica.carlsson@univ-lille.fr

Published in the Journal of Micropalaeontology, 41, 165–182, 2022 <u>https://doi.org/10.5194/jm-41-165-2022</u>

Version presented here: Published (Open Access)

Author contributions

Veronica Carlsson is the main author and writer of this paper. She formulated the research methodology in discussion with some of the other authors. Additionally, she classified, analyzed, collected, and processed all images and built up the datasets.

Taniel Danelian contributed to the discussion of formulating the research methodology and confirmed the classifications made by the main author. He also provided supervision, guidance, and reviewed the paper.

Pierre Boulet was involved in the discussion of formulating the research methodology and contributed to supervision, guidance, and reviewed the paper.

Philippe Devienne also participated in the discussion of formulating the research methodology and provided supervision, guidance, and reviewed the paper.

Aurelien Laforge contributed to this paper by coding the neural networks and the rotations of the images.

Johan Renaudie contributed to the discussion and reviewed the paper.

Abstract

This study evaluates the application of artificial intelligence (AI) to the automatic classification of radiolarians and uses as an example eight distinct morphospecies of the Eocene radiolarian genus *Podocyrtis*, which are part of three different evolutionary lineages and are useful in biostratigraphy. The samples used in this study were recovered from the equatorial Atlantic (ODP Leg 207) and were supplemented with some samples coming from the North Atlantic and Indian Oceans. To create an automatic classification tool, numerous images of the

investigated species were needed to train a MobileNet convolutional neural network entirely coded in Python. Three different datasets were obtained. The first one consists of a mixture of broken and complete specimens, some of which sometimes appear blurry. The second and third datasets were leveled down into two further steps, which excludes broken and blurry specimens while increasing the quality. The convolutional neural network randomly selected 85 % of all specimens for training, while the remaining 15 % were used for validation. The MobileNet architecture had an overall accuracy of about 91 % for all datasets. Three predicational models were thereafter created, which had been trained on each dataset and worked well for classification of Podocyrtis coming from the Indian Ocean (Madingley Rise, ODP Leg 115, Hole 711A) and the western North Atlantic Ocean (New Jersey slope, DSDP Leg 95, Hole 612 and Blake Nose, ODP Leg 171B, Hole 1051A). These samples also provided clearer images since they were mounted with Canada balsam rather than Norland epoxy. In spite of some morphological differences encountered in different parts of the world's oceans and differences in image quality, most species could be correctly classified or at least classified with a neighboring species along a lineage. Classification improved slightly for some species by cropping and/or removing background particles of images which did not segment properly in the image processing. However, depending on cropping or background removal, the best result came from the predictive model trained on the normal stacked dataset consisting of a mixture of broken and complete specimens.

2.1 Introduction

Polycystine radiolarians belong to an extant group of marine zooplankton protists secreting an aesthetically pleasing siliceous test that is rather well preserved in the fossil record and is therefore of importance to both biostratigraphy and paleoceanography. They are unique

amongst skeleton-bearing planktonic representatives in having a fossil record stretching as far back as the early Cambrian (Obut and Iwata, 2000; Pouille et al., 2011; Aitchison et al., 2017). Their continuous Cenozoic fossil record has allowed description of a number of welldocumented evolutionary lineages (Sanfilippo and Riedel, 1990), although their taxonomy has still not been fully clarified in spite of the great progress achieved during the last few decades (O'Dogherty et al., 2021). Polycystine radiolarian classification at the species level is based on morphological criteria, which therefore bear particular significance if one wishes to address evolutionary questions, but also for the development of high-resolution biostratigraphy. Supervised learning uses labeled data to train algorithms that will enable automatic classification and computer vision to deal with information from a visual context such as digital images or videos. These are some of the branches of artificial intelligence (AI) that have been developed during the last few years and may provide solutions to a number of difficult classification tasks. As such, convolutional neural networks (CNNs) use a deep learning algorithm to recognize patterns in images in a grid-like arrangement with multiple layers (Hijazi et al., 2015), which is a common approach for the analysis and classification of images. Training CNNs in a supervised way requires both a labeled training and validation dataset, from which the CNNs will learn features and patterns unique to each class from the training set by forming outputs, with which the untrained validation data will respond to if the model is a good fit.

A number of studies have attempted to apply automatic classification techniques on microfossils and/or micro-remains by using supervised machine learning in the past. Dollfus and Beaufort (1999) were the first micropaleontologists to apply AI in classifying and counting coccolithophores. They created the software SYRACO as an automatic recognition system of coccoliths, which was further developed a few years later (Beaufort and Dollfus, 2004) to count

automatically identified coccoliths but also for application to late Pleistocene reconstructions of oceanic primary productivity (Beaufort et al., 2001). Goncalves et al. (2016) tried to automatically classify modern pollen from the Brazilian savannah using different algorithms and achieved a highest median accuracy of 66 %, which is nearly as high as the median accuracy obtained by humans, based on a dataset that consisted of 805 specimens and 23 classes of pollen types. Hsiang et al. (2019) trained a neural network of 34 different modern species of foraminifera using a large dataset of a few thousand images, which reached an accuracy exceeding 87.4 %. Carvalho et al. (2020) used 3D images of 14 species of foraminifera, obtaining a dataset as large as 4600 specimens and a microfossil identification and segmentation accuracy as high as 98 % by using the CNN architectures of Resnet34 and Resnet50 with adjustment of hyperparameter optimization. De Lima et al. (2020) used a relatively small dataset of fusulinids composed of 342 images (including training, validation and test sets), which were divided into eight classes on a genus level to train in different CNNs architectures. They obtained the highest accuracy of 89 % by using the finetuning InceptionV3 model. Marchant et al. (2020) managed to train a CNN on a very large dataset of 13 001 images, including 35 different species of foraminifera. The best accuracy they obtained was about 90 %. Tetard et al. (2020) developed an automated method for new slide preparations, image capturing, acquisition and identification of radiolarians with the help of a software known as ParticleTrieur. They attempted to classify all common radiolarians existing since the Miocene, in a total of 132 classes, with 100 of them being relatively common species. They obtained an overall accuracy of about 90 %. Itaki et al. (2020) developed an automatization for the acquisition and deep learning of a single radiolarian species, Cycladophora davisiana, and obtained an accuracy similar to a human expert. Interestingly, they managed to be three times faster than a human being. Finally, Renaudie et al. (2018) applied the computationally efficient MobileNet convolutional neural network architecture for automatic radiolarian classification of 16 closely related species of the Cenozoic genera *Antarctissa* and *Cycladophora*. They obtained an overall accuracy of about 73 %, which they managed to increase to ca. 90 % after ignoring specimens which were not classified at all and by only including those specimens which had been given a class by the CNN.

The objective of our study was to obtain an accurate system of automatic classification for an automated classification of eight closely related species belonging to the middle Eocene genus *Podocyrtis* Ehrenberg, 1846 to be used by non-specialists in radiolarian taxonomy (e.g., students, industrial biostratigraphers or geochemists). Several of these species have a very good fossil record and are important in biostratigraphy as well as in morphometrics and evolutionary studies including gradual evolutionary transitions (Sanfilippo and Riedel, 1990; Danelian and MacLeod, 2019). In this work, we wished to implement MobileNet version 1 (Howard et al., 2017) because of its simplicity and lightweight construction, which enabled us to run more data in a shorter time. Many examples of MobileNet are available online and it is relatively easy to reproduce this work, which could thus be seen as a starting point for any other type of CNN implementation.

We will therefore attempt to answer the following scientific questions:

- 1. How well can the MobileNet convolutional neural network classify closely related species of the genus *Podocyrtis*?
- 2. How well can the predictive model classify *Podocyrtis* species under different processing settings and with materials coming from different parts of the world's ocean?

2.2 Materials

The main radiolarian material used for this study comes from the South American margin off Surinam (Demerara Rise, Ocean Drilling Program (ODP) Leg 207, Shipboard Scientific Party, 2004, see Table 1.1), where the middle Eocene interval is composed of an expanded sequence of chalk rich in abundant and well-preserved siliceous microfossils (for more details see Danelian et al., 2005, 2007; Renaudie et al., 2010). We focused on eight closely related species of the genus *Podocyrtis* (Fig. 2.1). Taxonomic concepts followed in this study are in accordance with Riedel and Sanfilippo (1970), Sanfilippo et al. (1985), Sanfilippo and Riedel (1990, 1992), with their stratigraphic ranges as specified recently by Meunier and Danelian (2022).

The eight studied species are considered to be part of three distinct evolutionary lineages, classified as three subgenera of the genus *Podocyrtis*: *Podocyrtis*, *Podocyrtoges* and *Lampterium* (Sanfilippo and Riedel, 1992). All taxonomic assignments were performed by a single taxonomist but were also checked by two other experts who had access to all 2D images of the dataset. The *Podocyrtis* species are in general relatively easy to recognize by their outer shape and/or size and distribution of pores (Figs. 2.1 and 2.2).



Figure 2.1. Age range and evolutionary relationships of *Podocyrtis* species occurring in Hole 1260A modified from Meunier and Danelian (2022). Arrows indicate descending species.

ODP Leg 207											
8		Radi									
		an	Р.	Р.	Р.	Р.	Р.	Р.	Р.	Р.	Specimens per
Hole	Sample	zone RP1	papalis	goetheana	chalara	mitra	sinuosa	ampla	phyxis	diamesa	sample
1259A	18R-1W, 53-55 cm	6	8	57							65
	18R-2W, 53-55 cm	554	16	56							72
	20R-3W, 53-55 cm	RP1 5	100		130	2					232
		RP1									
	21R-CC, 63-177 cm	and									
	and 25R-CC, 63-177 cm (mixed)	RP1 3	57			90	1				148
	26P 3W 54 56 cm	RP1	15				25			12	82
	20 R -5 W , 54-50 cm	RP1	15				25			42	02
1260A	6R-1W, 55-57 cm	6	2	4							6
	6R-2W, 55-57 cm			18	1						19
	6R-4W, 55-57 cm	RP1			2						2
	6R-CC, 63-177 cm	6/15 PP1	5		15						20
	7R-6W, 54-56 cm	5	4		2						6
	7R-CC, 63-177 cm				5	1					6
	8R-3W, 54-56 cm		7		20						27
	8R-5W, 54-56 cm				2						2
	8R-6W, 54-56 cm	DD1	1								1
	9R-1W, 55-57 cm	5/14	3								3
	9R-2W, 55-57 cm		11			14					25
	10R-5W, 55-57 cm	RP1 4	19			82		44			145
	13R-1W, 54-55 cm	RP1 4/13	10					2			12
	15R-4W, 55-57 cm	RP1 3/12	8				23		7		38
	16R-1W, 55-57 cm		5				26		57	17	105
	17R-1W, 55-57 cm		4							9	13
	17R-CC, 63-177 cm		10				10			26	46
	19R-6W, 55-57 cm						7				7
	22R-CC, 63-177 cm		3								3
Specime ns per species			288	135	177	189	92	46	64	94	1084

Table 2.1. Number of specimens per species for each sample coming from ODP Leg 207 which were used for training and validation.

The *Podocyrtis* subgenus represents an ancestral lineage that experienced morphological stasis. It is represented by the single morphospecies *Podocyrtis papalis* Ehrenberg, 1847, which differs from all the other *Podocyrtis* species by its partly developed abdomen (often smaller than the thorax), its overall fusiform shape (largest test width located on its thorax) and weakly expressed lumbar stricture. Three shovel shaped feet are often present, as well as a welldeveloped apical horn (which may be broken sometimes).

The *Podocyrtoges* subgenus is composed of three distinct morphospecies that belong to a lineage that evolved anagenetically. These are, from oldest to youngest:

Podocyrtis diamesa Riedel and Sanfilippo, 1970 differs from *P. papalis* by a more distinct lumbar stricture, with rather equally sized thorax and abdomen, and a more elongated than fusiform body. Some of the stratigraphically late forms of *P. papalis* display a degree of similarity in shape to *P. diamesa* (Fig. 2.2, 13th image), although the latter is much bigger in size (Fig. 2.1) and bears a larger apical horn than *P. papalis*.

Podocyrtis phyxis Sanfilippo and Riedel, 1973 displays a very distinct lumbar stricture formed at the junction between the abdomen and the thorax, with the former being more inflated than the latter. The overall outline of the test recalls the number eight "8". In general, a large horn is present on the cephalis. Complete specimens were rare in our material, as their horn is fragile and often broken.

Podocyrtis ampla Ehrenberg, 1874 displays a conical outline and a less prominent lumbar stricture than *P. phyxis*, with its abdomen being widest distally. Stratigraphically late forms of *P. ampla* do not display any feet and these forms were selected for this study.

The *Lampterium* subgenus is composed of four distinct morphospecies that belong to a lineage that also evolved anagenetically. These are, from oldest to youngest:

Podocyrtis sinuosa Ehrenberg, 1874 displays a barrel-shaped abdomen that is larger and more inflated than its thorax. Its widest part is located centrally at the mid-height of the abdomen.

Podocyrtis mitra Ehrenberg, 1854 displays an abdomen that is widest distally, rather than at mid-height as in *P. sinuosa*. It also displays more than 13 pores in the circumference of the widest part of the abdomen. Specimens with a rough surface on the abdomen, which could possibly be assigned to *P.* trachodes in the sense of Riedel and Sanfilippo, 1970, were included under *P. mitra*.

Podocyrtis chalara Riedel and Sanfilippo, 1970 displays a thick-walled abdomen, with large and more regularly arranged pores than *P. mitra*. It differs from the latter by having less than 13 pores in the circumference of the widest part of its abdomen. The *P. chalara* specimens selected for our material display 8 to 10 pores in circumference, so that they could be clearly distinguished from *P. mitra*.

Podocyrtis goetheana (Haeckel, 1887) displays long straight bars on its abdomen that enclose exceptionally large pores. The largest of them are located at the middle row of pores. They are often elongated, with four pores in the circumference. There are, however, some noticeably short specimens in our material that clearly belong to *P. goetheana* (Fig. 2.2, 4th image). No feet are present.

A total of 1085 radiolarian specimens were selected from the material available from the Demerara Rise. Their images were taken and prepared at the University of Lille and used for both training and validation of the CNN. The number of specimens used per species varies between 46 and 288 (Table 2.1). A second testing set of samples was prepared with 22 specimens that are stored at the Museum für Naturkunde in Berlin (see Table 2.2). Ten of them

come from the Indian Ocean (Madingley Rise, ODP Leg 115, Hole 711A, Shipboard Scientific Party, 1988), six other specimens come from the western North Atlantic Ocean (New Jersey slope, Deep Sea Drilling Project (DSDP) Leg 95, Hole 612, Shipboard Scientific Party, 1987) and six others also from the western North Atlantic Ocean (Blake Nose, ODP Leg 171B, Hole 1051A, Shipboard Scientific Party, 1998).

Table 2.2. Number of specimens per species for each sample coming from ODP Leg 115, DSDP Leg 95 and ODP Leg 171,which were used for test.

	Sample	Radiolarian zone	Podocyrtis specimens
ODP Leg 115			
Hole 711A	20X-2, 137-143 cm	RP16-RP15	5
	22X-2, 75-81 cm	RP15-RP14	3
	23X-1, 84-90 cm	\leq RP14	2
DSDP Leg 95			
Hole 612	20-1, 70-78 cm	RP16-RP14	2
	20-5, 106-116 cm		1
	22-6, 41-51 cm	\leq RP14	1
	27-1, 46-55 cm		1
	33-1, 45-56 cm		1
ODP Leg 171			
Hole 1051A	9H-04, 50-55 cm	RP16-RP15	2
	10H-02, 44-50 cm	\leq RP16	1
	31X-02, 44-50 cm	\leq RP14	1
	38X-02, 50-55 cm		2
Total number of specimens			22


Figure 2.2. A selection of the variety of *Podocyrtis* morphotypes analyzed in this study. 1-4) *P. goetheana*. 1-2) from 207_1259A_18R_1W. 53-55 cm; 3) 207_1260A_6R_2W. 55-57 cm; 4) 207_1259A_18R_2W. 53-55 cm. 5-6) *P. chalara* from 207_1259A_20R_3W. 53-55 cm. 7-9) *P. mitra* from 207_1260A_10R_5W. 55-57 cm. 10-12) *P. sinuosa*. 10) from 207_1259A_26R_3W. 54-56 cm; 11) 207_1260A_19R_6W. 55-57 cm; 12) Indian Ocean. 115_711A_25X_1. 83-86 cm. 13-14) *P. papalis*. 13-14) late (13) and typical (14) forms from 207_1259A_20R_3W. 53-55 cm. 15-16) *P. diamesa* from 207_1259A_26R_3W. 54-56 cm. 17-18) *P. phyxis* from 207_1260A_16R_1W. 55-57 cm. 19-20) *P. ampla* from 207_1260A_10R_5W. 55-57 cm. 55-57 cm. 19-20) *P. ampla* from 207_1260A_10R_5W. 55-57 cm. 55-57 cm. 19-20) *P. ampla* from 207_1260A_10R_5W. 55-57 cm. 55-57 cm. 19-20) *P. ampla* from 207_1260A_10R_5W. 55-57 cm. 55-57 cm. 19-20) *P. ampla* from 207_1260A_10R_5W. 55-57 cm. 55-57 cm. 19-20) *P. ampla* from 207_1260A_10R_5W. 55-57 cm. 55-57 cm. 19-20) *P. ampla* from 207_1260A_10R_5W. 55-57 cm. 55-57 cm. 19-20) *P. ampla* from 207_1260A_10R_5W. 55-57 cm. 5

2.3 Methods

We followed two different approaches for collecting photographs of specimens of *Podocyrtis* with the aim of constructing image datasets. A first approach involved the use of radiolarian slides from Leg 207, Hole 1260A, prepared initially for a biostratigraphic examination. The second involved the collection of *Podocyrtis* specimens, picked up directly and individually from dried residues of washed samples coming from both Holes 1260A and 1259A. The challenge faced while taking images of *Podocyrtis* from the old slides consists in specimens often touching themselves or overlapping with other objects. This led as a consequence to individual *Podocyrtis* specimens not being segmented properly by the methods described below and requiring manual (and time-consuming) segmentation.

2.3.1 Manual picking of individual Podocyrtis specimens

The residues used for sample preparation had already undergone acidic cleaning and removal of other non-siliceous particles by first dissolving the samples in hydrogen peroxide and afterwards in hydrochloric acid, followed by sieving at 50 μ m. Some of the samples needed further sieving to remove particles smaller than 45 μ m. They were then dried in a 50–60 °C oven.

Specimens of the various species of *Podocyrtis* were manually picked one by one under a ZEISS SteREO Discovery V20 microscope. The radiolarians were then transferred to a 32×24 mm coverslip and placed in such a way so as to avoid them being in contact with each other. A few drops of distilled water were placed on the coverslip for radiolarians to attach to the coverslip. Thereafter, they were dried overnight in an approximately 50 °C oven and thereafter attached on slides with Norland epoxy.

2.3.2 Image acquisition and processing

Images of *Podocyrtis* were taken under a Zeiss Axio A2 transmitted light microscope using the Zen 3.2 software with ×100 magnification and a pixel size set to 0.35 µm per pixel. Images were taken in fields of view (FOV), enabling several radiolarians to be captured at once in the same FOV. Approximately 3–15 focal points were taken on each FOV, depending on the specimen size. The images were then stacked using Helicon Focus 7 (Fig. 2.3a–c). They were then segmented (i.e., isolated from the background) all at once (Fig. 2.3c–d) using the ImageJ BioVoxxel plugin (Brocher, 2022) and a modified version of the Autoradio_Segmenter plugin (Tetard et al., 2020).

The images were then further processed with a script from Scikit Image version 0.18.1 (Van der Walt et al., 2014), using Python version 3.7.10 (Van Rossum and Drake, 2009), which rotated them in the same direction along a diagonal angle by finding the longest axis on the radiolarian specimen without cutting off objects within the picture (Fig. 2.3d–h). Thereafter all images were resized to equal 256×256 pixels. Having all images in the same orientation decreases the variability and increases the accuracy of the neural network. By having the images rotated in a diagonal angle optimizes the pixel resolution.

The most time-consuming task is the collection of images, as numerous images are needed for each species to build up a consistent dataset. Automatization of this task may be facilitated by the use of an automatic microscope, as in Itaki et al. (2020), Tetard et al. (2020) and Marchant et al. (2020).

The time needed for picking, slide mounting, photographing and image processing of around 100–200 radiolarians was one day. Manual picking speeds up the process, although caution

should be exercised to add glycerin or gelatin instead of water while mounting individually picked radiolarians on coverslips, in order to avoid formation of bubbles.



Figure 2.3. Image processing. (a) Three "raw" images taken from the microscope from the same FOV but different focal points, which are (b) stacked together using the Helicon Focus 7 software (c) to produce one entire focused and crisp image in which each particle in one FOV is segmented into an individual images or so-called vignette. (d) The segmented images are also transformed into a square 8-bit grayscale with a black background and white objects that are (e) further processed in Python by first rotating the radiolarian objects in the vignettes with a 45° angle so that the longest axis goes from the upper left corner to the lower right corner. (f) The vignettes are then filled again into squares so that no parts of the specimens are removed. (g)

Thereafter, they are cut again into squares just precisely so that the specimen fills the entire square and, lastly the new images are resized to 256 pixels on each side image (**h**).

2.3.3 Datasets

The radiolarian specimens included in the training and validation datasets contain only individuals that display a variability that may be included in the morphological boundaries accepted in the concept of each one of the eight species. Specimens that could not be identified with certainty as one or the other morphospecies (i.e., intermediate forms) were removed. All images are in full focus or so called "stacked". The images were divided into three different leveled datasets: the "normal" stacked dataset, called the "S" dataset (Fig. 2.4); the "SC" dataset, with only complete (unbroken) specimens, which obviously contains fewer radiolarian specimens but images of good quality; and the "SCUB" dataset, for which all blurry images were removed from the "SC" dataset (Fig. 2.4).

In all these datasets, 85 % of all specimens were used for training the model and 15 % of all specimens were only used to validate the trained neural network with the train and test split function from Scikit-learn (Pedregosa et al., 2011). This distribution aims to keep enough images to have quality learning, while having enough images for the network assessment to make sense of and to avoid miscalculations by running the model several times. Since the training and validation sets were randomized each time, it was important to perform several runs and thereafter take note of an average accuracy value. Neural network performance was compared between these different datasets.



Figure 2.4. Images of fully processed specimens (vignettes) with stacking, segmentation, rotation and resizing. (1) *P. ampla*.
(2) *P. chalara*. (3) *P. goetheana*. (4) *P. papalis*. (5) *P. diamesa*. (6) *P. phyxis*. (7) *P. mitra* (P. trachodes) and (8) *P. sinuosa*.
"a" and "b" stand for different individuals. S = stacked dataset. SC = stacked dataset with only complete specimens and SCUB = stacked dataset with only complete and unblurry specimens.

2.3.4 MobileNet convolutional neural networks

The CNNs are constructed by node layers including input layers which transfer their information in the form of node connections with different weights and threshold values into hidden layers, the convolutional layers, which process and transform the information into the next layers. Each convolutional layer has a different size and number of filters. A filter can be seen as a small grid of pixels, with different pixel values in each grid corresponding to a specific color value. This grid will go through an entire image in a sliding (convolving) way and transform the new values to the next layer that will process the image in a different or similar way. Early layers could, for example, easily detect edges, circles or corners, and later layers can even recognize more specific objects (Krizhevsky et al., 2012).

The network used for training specimens in classification is the Keras implementation of MobileNet (Howard et al., 2017), which is a convolutional neural network architecture (Chollet, 2015). Since MobileNet is a relatively small model, less regularization and data augmentation procedures are needed because smaller models have fewer problems with overfitting (Howard et al., 2017).

The input size of the images in the network was set to $256 \times 256 \times 3$ (number three stands for red, green and blue (RGB) colors or channels), dropout was set to 0.15 for layers trained by ImageNet inside the MobileNet architecture and an average pooling was used. An added dropout layer was set to 0.5 and added after the MobileNet convolutional layers to prevent overfitting by randomly switching off some percentages of neurons in the model. Finally, a dense layer was added, which is the most commonly used layer in neural network models. It performs a matrix-vector multiplication, for which values are parameters and which can be trained and updated with backpropagation, and the dense layer was set to eight outputs corresponding to the number of species. The SoftMax activation used here converts the values into probabilities. The optimizer used was "Adam", a stochastic gradient descent, and the loss function used was "Categorical Crossentropy". The batch size of 64 resulted in 100 steps per epoch and only three epochs were necessary for the training, for the simple reason of avoiding any overfitting models. After three and sometimes four epochs, the validation accuracy does not increase further, and the loss becomes bigger (see Tables 2.S1-2.S3 in the Supplement 2.1 for an example of a MobileNet run on five epochs). For each dataset, since we used the traintest split function, the model was run 10 times to obtain good enough average accuracies.

2.4 Results

2.4.1 CNN accuracies

The MobileNet convolutional neural network model used here resulted in an average validation accuracy of 88.46 % for 10 runs of the "S" dataset, 92.13 % average accuracy for the "SC" and 92.39 % for the "SCUB" datasets for 10 runs on each dataset (Table 2.3). It is important to investigate how each run was performed, since it may vary a lot, especially by looking at each individual species' performance (Table 2.3). Although there are codes that can equally select 15 % from each species, we chose not to use that option here because we also wanted to see how the model performs without selecting all general forms for each species. The total time to run MobileNet on all datasets (total of 30 times) was around 8.5 h, 10–15 min for each run. In general, the "S" dataset with complete specimens had the smallest difference between the lowest and highest accuracies over its 10 runs. The dataset consisting of stacked, complete, clear and crisp specimens also had low variation between the highest and lowest accuracies.

 Table 2.3. Validation accuracies for all the studied species over 10 runs for each one of the three datasets, and their average values based on results of the species accuracies from Supplement 2.2.

S									
Run	P. ampla	P. chalara	P. diamesa	P. goetheana	P. mitra	P. papalis	P. phyxis	P. sinuosa	Average
1	100.00	100.00	33.33	100.00	90.91	88.89	100.00	92.86	87.88
2	100.00	67.74	100.00	94.74	100.00	95.35	80.00	80.00	89.69
3	100.00	100.00	27.78	100.00	78.57	31.25	57.14	53.33	62.03
4	100.00	95.65	64.29	100.00	100.00	53.66	87.50	100.00	83.54
5	100.00	100.00	93.33	100.00	100.00	92.86	100.00	100.00	97.47
6	100.00	100.00	100.00	95.65	100.00	88.37	100.00	100.00	96.38
7	100.00	100.00	44.44	90.48	100.00	97.96	77.78	100.00	91.50
8	100.00	100.00	78.57	100.00	100.00	97.30	93.33	93.33	96.38
9	100.00	100.00	100.00	85.71	62.50	75.00	71.43	88.89	83.91
10	100.00	100.00	0.75	91.67	96.43	96.15	100.00	100.00	95.84
Average	100.00	96.34	64.25	95.82	92.84	81.68	86.72	90.84	88.46

SC									
Run	P. ampla	P. chalara	P. diamesa	P. goetheana	P. mitra	P. papalis	P. phyxis	P. sinuosa	Average
1	100.00	100.00	80.00	88.89	100.00	100.00	100.00	100.00	96.46
2	100.00	100.00	100.00	100.00	100.00	87.10	100.00	92.31	94.69
3		100.00	75.00	100.00	100.00	94.12	100.00	100.00	96.46
4	100.00	100.00	72.73	100.00	100.00	100.00	100.00	100.00	97.35
5	100.00	100.00	100.00	100.00	100.00	96.97	100.00	100.00	99.12
6	100.00	100.00	83.33	100.00	100.00	100.00	0.00	0.00	99.12
7	75.00	70.00	0.00	100.00	12.50	60.61	50.00	100.00	47.79
8	100.00	100.00	100.00	100.00	100.00	79.50	0.00	100.00	99.12
9	100.00	100.00	66.67	100.00	93.10	80.77	100.00	100.00	92.04
10	100.00	100.00	90.00	100.00	100.00	100.00	100.00	100.00	99.12
Average	97.22	97.00	76.77	98.89	90.56	89.91	75.00	89.23	92.13
SCUB									
Run	P. ampla	P. chalara	P. diamesa	P. goetheana	P. mitra	P. papalis	P. phyxis	P. sinuosa	Average
1	100.00	100.00	100.00	100.00	100.00	84.62	100.00	50.00	94.37
2	80.00	75.00	60.00	100.00	100.00	70.00	100.00	100.00	85.92
3	100.00	100.00	100.00	100.00	88.24	100.00	100.00	100.00	97.18
4		100.00	100.00	100.00	100.00	93.33	100.00	85.71	97.18
5	100.00	100.00	100.00	100.00	100.00	94.12	94.12	100.00	98.59
6	100.00	100.00	80.00	100.00	100.00	100.00	100.00	100.00	98.59
7	100.00	100.00	71.43	100.00	81.82	81.82 81.82 0.00		100.00	92.96
8	100.00	100.00	0.00	100.00	57.14	57.14 77.78 33.33 40.0		40.00	73.24
9	100.00	94.12	62.50	100.00	100.00	100.00		100.00	94.37
10	100.00	100.00	75.00	100.00	93.75	93.75 88.24 100.0		75.00	91.55

2.4.2 Confusion matrices

The neural networks also produced confusion matrices for each run. Since the training and validation sets were randomized, we therefore created three average confusion matrices (Fig. 2.5), one for each type of dataset. Since the number of specimens used for validation varies, we calculated an average value for the validation size as well. In these confusion matrices the y-axis shows the actual species, while the x-axis shows predicted ones. Each box shows the average accuracy based on the validation set.

What is further observed is the fact that closely related species (i.e., morphospecies situated along an evolutionary lineage) are often mistaken for each other. This is true especially for those species with more than one neighboring species along a lineage, which is the case for all species studied here with the exception of the lineage end members, e.g., *P. ampla* and *P. goetheana*. A very remarkable point is that *P. diamesa* appears to often be misidentified as *P. papalis*, more frequently than *P. papalis* is misinterpreted as *P. diamesa*, which results in the average precision value being lower in *P. diamesa* compared to the rest of the species, while *P. papalis* has a lower recall than precision value.



Figure 2.5. Average confusion matrices for the "S", "SC" and "SCUB" dataset. The numbers inside the matrices shows the average validation accuracy of specimens in each class that has been correctly or incorrectly classified. The numbers under each species name are the total number of true or predicted labels.

Figure 2.6 displays the calculated average precision, recall and F1 score for all species based on the confusion matrices. The precision value also means the correct prediction value, and can be simplified by the following Eq. (1):

$$Precision = \frac{\text{Number of specimens classified as a class and also belonging to that class}}{\text{All specimens classified as that particular class}}$$
(1)

The recall values show that not all specimens belonging to a class have been classified to the correct class, similar to the accuracy, which is the number of specimens correctly classified divided by the total number of specimens, as in the Eq. (2):

$$Recall = \frac{\text{Number of specimens in a class that were correctly classified}}{\text{Total number of specimens in that class}}$$
(2)

The F1 score is an average value of this and shows the average between precision and recall written like Eq. (3):

$$F1 = \frac{2*(\text{Precision}*\text{Recall})}{\text{Precision}+\text{Recall}}$$
(3)

The average F1 scores based on the confusion matrices are 89.05 %, 93.26 % and 91.72 % for the "S", "SC" and "SCUB" datasets, respectively. This implies that the best result for the F1 score is obtained when the CNN was trained on the "SC" dataset.



Average Precision, Recall and F1 score

Average Precision, Recall and F1 score SC dataset





Figure 2.6. Average precision, recall and F1 score calculated from the average confusion matrices based on the S, SC and SCUB datasets.

2.4.3 Testing the predictive models

Once a high validation accuracy (i.e., over 90 %) had been obtained for each dataset following the training process, parameters were saved and predictive models were formed. The idea is to explore how the classification tool classifies Podocyrtis species and with what level of confidence. Therefore, a total of 22 specimens of *Podocyrtis* obtained from ODP Leg 171B, Hole 1051A (Blake Nose, western North Atlantic Ocean), DSDP Leg 95, Hole 612 (New Jersey slope, North Atlantic Ocean) and ODP Leg 115, Hole 711A (Madingley Rise, Indian Ocean) were used as a test dataset. Slides were prepared with Canada balsam and photographs were taken with a Leica transmitted light microscope to which an AmScope MU1003 digital camera was attached. Images were then segmented with ImageJ, and rotated and resized with a Python script as described above. Some images were further cropped, but no background particles were removed. Tables S4-S6 in Supplement 2.1 display how each specimen from the different locations was classified using different predictive values, a comparison with classification by cropping images that appeared very tiny in relation to the entire image and a comparison of removing background particles of those specimens which displayed them. The main result shows that P. sinuosa is often misinterpreted as P. papalis even though P. papalis is almost always classified correctly. There are significant morphological differences between the species trained in the neural network and the test set. Here, the best predictive model to use is the "S" dataset and the worst one is "SCUB" (Table 2.4), which is quite the reverse from the validation set based on the tropical Atlantic specimens from the Demerara Rise trained in this network.

Results obtained after the model "S" was applied on material from the North Atlantic and Indian Oceans were 13 out of 22 specimens classified correctly without any manual fixation, which corresponds to 59 % accuracy. The accuracy of model "SC" was raised to 68 % after manually

cropping images for specimens that appeared smaller, while it was increased up to 77 % by removing background particles appearing in the images. The "SC" model produced worse results. For all images without any manual fixation, the accuracy obtained was 45 %, but increased up to 50 % after adding manual cropping and to 54.5 % after adding the segmentation. The "SCUB" model had an accuracy of 41 % for all three of the different image fixations. However, in all cases, in terms of neighboring species, at least 20 specimens were correctly classified as a neighboring species which translated into an accuracy of at least 90.9 %.

 Table 2.4. Result of all the 22 specimens including necessary manual cropping and segmentation from ODP Leg 115, Hole

 711A from the Indian Ocean; DSDP Leg 95, Hole 612; and ODP Leg 171, Hole 1051A from the western North Atlantic

 Ocean which were classified using different parameters trained on the three different datasets, "S", "SC" and "SCUB".

									Manual			
	All			Manual cropping					segmentation			
				Correct			Correct					Correct
Classification	Correct	Incorrect	Uncertain	neighbouring	Correct	Incorrect	Uncertain	neighbouring	Correct	Incorrect	Uncertain	neighbouring
				species				species				species
S	13	7	2	20	15	5	2	21	17	3	2	21
SC	10	10	2	21	11	9	2	21	12	8	2	21
SCUB	9	11	2	22	9	11	2	22	9	11	2	21
Total number	22											
of specimens	22											

2.5 Discussion

2.5.1 MobileNet performance and accuracy

Dedicated to embedded systems and smartphones for which low latency and real-time execution are key demands, the advantage of using the MobileNet architecture is that it is extremely light and small (in terms of coding and weight of models). It is fast, with an only slight degradation in inference accuracy according to the gain of the consumed resources, and easily configurable to improve detection accuracy (Howard et al., 2017). When tested on Im2GPS, a dataset which

gives geolocation from images taken from different spots around the world, the accuracy of MobileNet was about 1 % higher compared to GoogleNet, whereas it used 2.5 times less computation and, as cited in Howard et al. (2017), "MobileNet is nearly as accurate as VGG16 while being 32 times smaller and 27 times less compute intensive". The study by Howard et al. (2017) presents extensive experiments on resource and accuracy trade-offs and shows strong performance of MobileNet compared to other popular models on ImageNet classification. This is the reason why some works (Rueckauer et al., 2021) start to deploy MobileNet also on neuromorphic hardware such as Loihi (Davies et al., 2021). Although the development of AI has been based until now on software bricks installed on big data centers, the current multiplication of connected objects requires decentralization. The new AI revolution now involves development of specific electronic components with very promising results.

The images that we have used were transformed to RGB-colored because the pretrained weights of ImageNet are only compatible for RGB-colored images, as this is also the case of the whole architecture of MobileNet; the idea here is to apply a depthwise convolution for a single filter for each unique input channel. The use of neural network models, equally adapted for grayscale images, could also decrease the energy consumption greatly. In terms of resolution, MobileNet resizes images into a lower resolution. In most cases this did not affect the result, but it is plausible that in a few cases the neural network was not able to see the position of the lumbar stricture, which is an important distinguishing character. In any case, Renaudie et al. (2018) also commented on the resolution loss due to resizing, as the inner spicules in *Antarctissa* species disappeared, which are crucial for species identification.

The species with the highest F1 score (Fig. 2.5) were *P. goetheana* (94 %–97 %) and *P. ampla* (90 %–99 %). A reason for this might be that both species are at the end of the *Lampterium* and

Podocyrtoges lineages and have only one closely related species, while the others have two. For P. chalara, P. sinuosa and P. mitra, a reduction in the number of analyzed specimens increased the precision and increased the recall value for *P. mitra*, giving *P. mitra* an overall better F1 score with reduction of specimens. This is expected, since the variability decreases when removing more imperfect specimens but performs better for determining unbroken and clear specimens. A reduction of specimens for P. goetheana did not make a big difference to precision but it did increase the recall. In general, P. diamesa has both the lowest recall (69 %-85%) and precision (80%-84%) values. Although P. diamesa and P. papalis are often mistaken for one another, it is mostly P. diamesa that is misidentified as P. papalis, which may be due to the fact that late morphotypes of P. papalis resemble P. diamesa to some extent. The distinguishing character of this species is that the overall shape of *P. papalis* is in most cases rather fusiform with a larger thorax than abdomen, which is only partially developed. Podocyrtis diamesa is much larger (even though this does not seem to matter since all images are resized to equal sizes), the size of the thorax and abdomen is often more similar, and the lumbar stricture is more prominent. Although late forms of P. papalis do resemble P. diamesa, they do not co-exist at the same time interval (the true P. diamesa morphospecies disappeared well before the appearance of the late *P. papalis* forms) and they are smaller.

The highest F1 score for all species except *P. diamesa* comes from using a model trained on complete specimens ("SC"), regardless of quality. *Podocyrtis phyxis* is one example of a species that has the best performance in the "SC" dataset. The overall dataset of *P. phyxis* consists of many specimens that are broken and missing the apical horn; therefore, there is a significant reduction of specimens when using the datasets with only complete specimens. The number of specimens present for *P. phyxis* in the "S" dataset was 64. This number dropped down to 22

specimens in the "SC" dataset and then further to 16 specimens in the "SCUB" dataset. When it comes to *P. sinuosa*, there is a substantial reduction due to many blurry specimens, which is likely the result of the mounting media. A total of 49 blurry specimens were removed from the "SC" dataset (78) compared to the "SCUB" dataset (29), and the overall result seems to have been improved somewhat by removing specimens.

To summarize, our work has about 91 % accuracy if we exclude the tests run with the unstacked (U) dataset (Supplement 2.1, Table S7), which not only consisted of blurry unfocused images but also kept some touching particles. Nevertheless, it has an overall similar accuracy of about 72 % to the work produced by Renaudie et al. (2018) with the same neural network. These authors used images as they appear under the microscope, without stacking or image processing. To produce more datasets with these types of images, one can expect a slightly lower accuracy than when images are processed. Renaudie et al. (2018) also included a substantial number of unidentifiable specimens. If these unidentifiable images were ignored, results would probably come close to 90 % accuracy. In our case, all specimens are identifiable by our CNN as any of the eight possible *Podocyrtis* species, even if the specimen in question is not a *Podocyrtis*. A solution for this issue is to perhaps apply parallel neural networks with a hierarchical architecture, similar to the one that Beaufort and Dollfus (2004) applied for SYRACO. One suggestion could be first to classify radiolarians and non-radiolarian particles, with a second step classifying radiolarians into higher taxonomic orders (Spumelleria, Nassellaria and unidentified broken radiolarians) and finally, a last step leveling down to species, genus and/or family levels. Itaki et al. (2020) went with another approach. These authors focused on the identification of one single species, Cycladophora davisiana, but they also used the morphologically similar species, Cycladophora bicornis, as another class, to avoid this species being wrongly interpreted as C. davisiana. Thereafter, they used the classes "centric diatoms", "all other radiolarians" and "all other objects".

It is worth noting that validation accuracies and F1 score values may appear high since we do not have a large dataset and the images in the datasets are in the same orientation and rotational angle, centralized in the middle. As mentioned earlier, species are carefully selected to avoid including any intermediate forms in the dataset. 5.2 Predictive models Given that our initial analysis was performed with material coming from the equatorial Atlantic (ODP Leg 207), we wished to consider a different dataset to test the predictive models, which consisted of images of specimens coming from the North Atlantic (DSDP Leg 95, ODP Leg 171) and the Indian Ocean (ODP Leg 115). In addition, radiolarians were mounted in a different mounting media (Canada balsam) and we used a different microscope. In most cases, particles were segmented properly in the segmentation process. However, some specimens were still attached to other particles, which could also result in that the specimens did not fill out the entire images. To save time and effort, we did not crop or remove background particles at first. These specimens could still be classified correctly. The results improved after the specimens were cropped or had their background particles removed, at least by using the predictive model trained by the "S" dataset. Podocyrtis papalis is almost always correctly interpreted. The reason might be the large number of specimens used in our dataset and the large number of morphological differences within the species. As mentioned earlier P. sinuosa is often misinterpreted as P. papalis. This is almost always the case for all P. sinuosa, whether we use the "SC" or "SCUB" datasets. The main reason for this probably lies in the fact that the morphotypes of P. sinuosa coming from the North Atlantic and the Indian oceans differ significantly (Fig. 2.2) from the ones trained in the initial neural network. Two other factors may contribute to this: first the decrease in *P. sinuosa* specimens when passing from the "S" dataset to the "SC" and "SCUB" datasets, as explained above; and secondly, the image quality, since *P. sinuosa* appear very whitish in the trained neural network. One *P. chalara* was first classified as *P. mitra* before it was cropped. One possible reason may have been that the pores of the uncropped version appeared smaller, as in *P. mitra*. After cropping out unnecessary space, *P. chalara* appeared larger and could be correctly classified. In one case, one specimen, which was clearly a late *P. mitra* (according to preference of the author and not *P. trachodes*), was completely wrongly identified as either *P. phyxis* by using the "S" predictive model, *P. ampla* by using the "SC" predictive model or *P. goetheana* by using the "SCUB" predictive model. The reason for this may be due to the mounting media or preservation, because the pore space appears blacker or cleaner, similar to *P. phyxis* or *P. ampla*, which are generally larger from the tropical Atlantic assemblages trained in this network, and resizing the images may make them appear to be in a better resolution, with no white "dirt" between the pores and within the pore space. *Podocyrtis goetheana* also have gigantic pores and a lot of black space.

Most often it is closely related species that are mistaken for each other, as observed in the training and validation. It can be observed in Table 2.4 that even if specimens were not always interpreted as the right species, they were almost always interpreted as a neighboring species along a lineage. It was also discussed by Renaudie et al. (2018) that closely related species tend to be misinterpreted as each other due to morphological similarities, which is also confirmed in this work.

2.5.2 Species choice and their image properties

Akin to the study of Renaudie et al. (2018) conducted on Neogene radiolarians, we used the MobileNet neural network to classify closely related species of the Eocene genus *Podocyrtis*.

We did not, however, use images as they appeared under the microscope, since we used software and codes for image processing which can easily process several images into equal settings at once and do not only increase the image quality but also save a lot of time and effort.

Renaudie et al. (2018) chose to select all specimens present in a slide that could somehow be classified with reliable confidence. The same approach was followed here, but most of the samples were pre-selected knowing that some typical morphotypes existed in them. We used specimens for which species identification and classification was certain in most cases, meaning specimens which could be instantly recognized to one species, and left out uncertain intermediate forms. We also used specimens which were preserved nearly completely, even though many experts are often able to identify specimens based on even small fragments to at least a genus level. Smaller identifiable fragments would probably require a huge amount of data to train. The samples obtained here have excellent preservation and finding broken fragments in, for example, dinoflagellate cysts seem more likely than finding broken pieces of radiolarians, provided they have not been crushed by mounting.

2.6 Conclusions

The goal of this study was to create an automatic classification tool to allow AI-based identification of middle Eocene *Podocyrtis* species which would achieve the highest possible accuracy after training the MobileNet CNN based on a dataset of 1085 images of *Podocyrtis* morphotypes classified as eight different species.

Regarding our first question stated in the introduction, "How well can the MobileNet convolutional neural network classify closely related species of the genus *Podocyrtis*?", we showed that specimens which belong to *Podocyrtis* species can be classified automatically with

a high accuracy (91 % confidence). Best results were obtained by using datasets with improved quality (but a smaller number of images), both according to overall accuracy and the F1 score values. However, tests on *Podocyrtis* species from the North Atlantic and Indian Ocean work best by using the predictive model trained by the normal stacked dataset, consisting of more specimens but with a mix of broken, complete, blurry and clear images. This suggests that a higher variance of morphotypes could be applied to the datasets. In conclusion, this identification tool works well for classification of *Podocyrtis* species, although it could still be further improved by adding additional closely resembling species of *Podocyrtis* that were not present or very rare in our material.

Regarding our second scientific question, "How well can the predictive model classify *Podocyrtis* species under different processing settings and using material coming from different parts of the world's oceans?", we establish that the predictive models also work well for classifying images taken by different microscopes, but might in some cases require adjustment of the clarity settings and images taken by different mounting medias.

This study could be further improved by including additional morphospecies of *Podocyrtis* in the datasets and more specimens and morphotypes, especially from many other different oceanic realms. Another improvement to the neural network would be to detect *Podocyrtis* species or other taxa of interest among hundreds to thousands of other objects. This could perhaps be solved by classifying every object, as, for example, done in Tetard et al. (2020) and Itaki et al. (2020), or by applying a parallel network approach (Beaufort and Dollfus, 2004) first that filters away objects of no interest in different steps. For example, as a first step, this would involve the classification of radiolarians and non-radiolarians, and as a second step, the exclusion of all non-radiolarians. It could also be beneficial to create a network inspired by

MobileNet but adapted to grayscale images, in order to become more energy efficient and have a more appropriate resolution that detects small crucial details in radiolaria classification.

Data availability

Microscope slides from Leg 207, Hole 1259A and 1260A, which were used for training and validation of the neural network are stored at the University of Lille, France, and slides from ODP Leg 115, Hole 711A, DSDP Leg 95, Hole 612 and ODP Leg 171B, Hole 1051A, which were used for testing the CNN and are stored at the Museum für Naturkunde in Berlin, Germany. Datasets (https://doi.org/10.57745/G7CHQL, Carlsson, 2022) and codes (https://doi.org/10.57745/J4YL4I, Carlsson and Laforge, 2022) are published in the University of Lille repository at Recherche Data Gouv.

Supplement

The supplement related to this article is available online at: <u>https://doi.org/10.5194/jm-41-165-</u> 2022-supplement.

Acknowledgements

This study was supported by the French government through the program "Investissements d'avenir" (I-ULNE SITE/ANR-16-IDEX-0004 ULNE) managed by the National Research Agency. This project received funding from the European Union's Horizon 2020 research and innovation program under the Marie Sklodowska-Curie grant agreement no. 847568. It was also supported by UMR 8198 Evo-Eco-Paléo and IRCICA (CNRS and Univ. Lille USR-3380). We would also like to give special thanks to Mathias Meunier for giving opinions on taxonomy

and thanks to Hammouda Elbez who was a huge help with coding issues. Also, a huge thanks to David Lazarus for discussion on the topic and for facilitating the access of VC to the Museum

für Naturkunde in Berlin to take images of material from ODP Leg 171B, Hole 1051A; DSDP Leg 95, Hole 612; and ODP Leg 115, Hole 711A.

<u>CHAPTER 3</u>: Morphometrics and machine learning discrimination of the middle Eocene radiolarian species *Podocyrtis chalara*, *Podocyrtis goetheana* and their morphological intermediates

Francisco Pinto ^a, Veronica Carlsson ^{a,b*}, Mathias Meunier ^a, Bert Van Bocxlaer ^a, Hammouda Elbez ^b, Marie Cueille ^a, Pierre Boulet ^b, Taniel Danelian ^a

^a Univ. Lille, CNRS, UMR 8198, Evo-Eco-Paleo, F-59000 Lille, France

^b Univ. Lille, CNRS, CRIStAL – Centre de Recherche en Informatique Signal et Automatique de Lille, UMR 9189, F-59000 Lille, France

* Corresponding author. E-mail address: veronica.carlsson@univ-lille.fr

Published in Marine Micropalaeontology, 185, 2023 https://doi.org/10.1016/j.marmicro.2023.102293

Version presented here: Accepted (No Open Access)

Author contributions

Francisco Pinto collected and classified the images and wrote the initial draft of the manuscript; he led the work on the description of the intermediate forms; he conducted the LDA analysis and participated in the neural network analysis.

Veronica Carlsson is the corresponding author and formulated the research methodology in discussion with the other authors. She also confirmed the classification of end member and intermediate species and contributed with more images, took part in the writing and reviewing of the manuscript; she led the discussions and analyses on neural networks.

Mathias Meunier contributed to the supervision of image collection and preparation and led the actions on their LDA analysis. He also contributed to the writing and revision of the manuscript.

Bert Van Bocxlaer contributed to the design of the research methodology, especially concerning the LDA analyses, and has been part of the analytical work and discussion. Has was also involved in the writing and revision of the manuscript.

Hammouda Elbez coded and performed the neural network experiments and he contributed to the writing, discussion and revision regarding the neural network part.

Pierre Boulet contributed to the discussion and supervision, especially regarding the neural network part.

Marie Cueille contributed to the initial steps of this study that led to the discovery and first description of the intermediate morphotypes.

Taniel Danelian initiated this study and led the discussions that formulated the research methodology; he supervised the whole project and contributed to the writing and revision of the manuscript.

Abstract

We present various approaches to distinguish the middle Eocene species *Podocyrtis chalara* and Podocyrtis goetheana, which are end members of a trajectory of phenotypic change, and their intermediate morphogroups. We constructed a set of thirteen traditional morphological variables to classify the entire morphological variability encompassed by the two morphospecies and their intermediates Podocyrtis sp. cf. P. chalara and Podocyrtis sp. cf. P. goetheana. We used two methods of classification, namely Linear Discriminant Analysis (LDA) and machine learning using artificial neural networks. LDA performed on the morphometric data reveals a good discrimination for P. chalara, P. goetheana and Podocyrtis sp. cf. P. goetheana, but not for Podocyrtis sp. cf. P. chalara. We used three approaches of machine learning based on different neural networks: Convolutional Neural Networks (CNNs) and two Spiking Neural Networks (SNNs). Each of these neural networks was trained based on classified images of the two morphospecies and their morphological intermediates, thus constituting a different set of input data than the morphometric dataset for LDA. The neural network approaches identified the same three morphospecies recognized by LDA from a dataset of traditional measurements, i.e. P. chalara, P. goetheana and Podocyrtis sp. cf. P. goetheana, with up to 92 % accuracy. Our results highlight the great potential and promising perspectives of machine learning and neural networks in the application of image-based object recognition for morphological classification, which may also contribute to more objective taxonomic decisions.

Keywords: Morphometrics; Artificial Intelligence; Convolutional Neural Networks; Spiking Neural Networks; Radiolarians; Automated identification

3.1 Introduction

Polycystine radiolaria are one of the oldest known Rhizarian lineages, with a fossil record stretching back to the Early Cambrian (Obut and Iwata, 2000; Pouille et al., 2011) and is thus of much interest for a number of evolutionary studies (e.g., Danelian and Johnson, 2001, Danelian et al., 2014, Renaudie and Lazarus, 2003, Tetard et al., 2017).

Since the early stages of the Deep-Sea Drilling Program, the Cenozoic record of polycystine radiolaria has allowed us to establish their evolutionary and biostratigraphic significance, especially based on representatives of the Eocene genus *Podocyrtis*. (Sanfilippo and Riedel, 1970; Riedel, 1971; Moore, 1972; Riedel and Sanfilippo, 1978). Calibrated initially to the magnetostratigraphic time scale (Sanfilippo and Nigrini, 1998), the middle Eocene tropical radiolarian zones are now tied to orbital chronology (Meunier and Danelian, 20222), which provides the highest resolution of temporal control possible today and allows to define biostratigraphic events more accurately. Indeed, many of the middle Eocene biozones are based on the evolution of the various lineages of the genus *Podocyrtis*, which often relate to gradual anagenetic changes in phenotypes, as documented in several evolutionary lineages since the early 1970s (Sanfilippo and Riedel 1970, 1992; Riedel, 1971). For example, the bases of biozones RP14 and RP15 are defined based on anagenetic phenotypic changes between the morphospecies *P. sinuosa - P. mitra* and *P. mitra - P. chalara*, respectively.

The *Podocyrtis (Lampterium)* lineage ends with the marked morphological transition of *Podocyrtis chalara* to *P. goetheana*. Interestingly, the first occurrence of *P. goetheana* defines

the base of biozone RP16 (Sanfilippo and Nigrini, 1998); an anagenetic transition was reported, but intermediate forms were never documented in detail. The absence of such a documentation has implications for our understanding of evolutionary changes in this *Podocyrtis* lineage, but as intermediate forms are poorly understood it also affects the recognition of the base of RP16. Here we examine the well-preserved radiolarian material of ODP Sites 1259 and 1260 from Demerara Rise (equatorial Atlantic Ocean), which present an exceptionally expanded Eocene sedimentary sequence. As such, this material offers an exceptional opportunity to study the morphological transitional forms between *P. chalara* and *P. goetheana*.

Within the abovementioned context, the principal aim of our study is to document morphological variation in the anagenetic sequence of *P. chalara* to *P. goetheana* with two different but complementary approaches to test the performance of various machine learning algorithms based on neural networks. To reach this objective, we first quantified morphological variation in the anagenetic transition between *P. chalara* and *P. goetheana* with traditional morphometrics, i.e., linear measurements, pore counts and associated ratios in the framework of qualitatively recognized morphological entities. This quantification of shape follows previous attempts to assess morphological changes in the *Podocyrtis (Lampterium)* lineage (Sanfilippo and Riedel, 1990; Rohlf and Bookstein, 1991; Danelian and Macleod, 2019; Watanabe et al., 2022). Using this morphometric framework of measurement data and *a priori* morphospecies assignments, we examined how well linear discriminant analysis allows to distinguish morphospecies as a baseline to test the performance of machine learning with neural networks.

Testing the capabilities of neural network approaches based on image recognition is a daunting task, because during the past couple of years, a variety of techniques involving Artificial Neural

Networks (ANNs) have been developed and improved. Convolutional Neural Networks (CNNs) have been specifically designed for analysis of visual data, and are now commonly used for image recognition, warranting detailed examinations of their performance in morphological classification, and, therewith, as a tool to inform, and potentially reach more objective taxonomic decisions. Indeed, CNNs are becoming well-integrated in various micropaleontological studies for automatic image recognition (i.e., Mitra et al., 2019, Hsiang et al., 2019; Marchant et al., 2020; Dollfus and Beaufort, 1999; Beaufort and Dollfus, 2004; Bourel et al., 2020; Itaki et al., 2020; Renaudie et al., 2018; Tetard et al., 2020). Regarding Eocene radiolaria, a recent study by Carlsson et al. (2022) applied a CNN to eight welldelimited morphospecies of the genus Podocyrtis, and documented the potential of this method under the simplified scenario when no morphological intermediates are present. Spiking Neural Networks (SNNs) present another type of neural network, which in addition to neuronal and synaptic states, they also incorporate a time component; this is why such networks can more closely mimic natural neural networks (Maass, 1997). SNNs have wide applicability, including modeling of natural systems such as the central nervous system of biological organisms, as well as for image analysis. Traditionally, SNNs were less accurate than other neural networks, but in recent years their performance has significantly improved; they are more appropriate to process spatio-temporal data, and they may use computational resources more effectively (Tavanaej et al., 2019). As such, an evaluation of SNNs in image recognition and biological classification seems warranted.

In this paper, we expand on previous work by Carlsson et al (2022) with CNNs by aiming to classify stacked and segmented images of the entire spectrum of morphological variation found between *P. chalara* and *P. goetheana*. As mentioned, we used morphometrics to document

shape variability, which we subjected together with *a priori* morphospecies assignments to LDA as a baseline to study the classification performance for image-based neural networks using a CNN, a Spike-timing-dependent plasticity (STDP)-based SNN and a SuperSpike-based SNN. This examination allows us to evaluate the use of imaging data and neural networks for automated classification in a complex case study with intermediate shapes. If the neural networks perform well, we would expect correct classification for each morphogroup. Additionally, the results should reflect those of the LDA analysis, if the morphometric documentation of shape variation is representative of the four morphospecies. Alternatively, neural networks may show differences compared to LDA. These distinctions could arise if neural networks fail to perform well, possibly due to unsuccessful training with the existing data. In opposite to that, the neural networks may perform better than LDA if the data supporting LDA lacks crucial shape information necessary for distinguishing between morphogroups based on the images. As such, we expect our study to shed light into future opportunities for automated biological classification of polycystine radiolaria and the use of neural networks in developing more objective taxonomic decisions.

3.2 Analyzed morphological groups

Plate 3.1 displays the entire range of morphological variability observed between *P. chalara* and *P. goetheana*. As linear discriminant analysis and the supervised learning of neural networks are based on *a priori* group assignments, we were required to assign this continuum of variation to a number of morphological groups. Based on extensive qualitative assessments and to challenge the employed classification algorithms we recognized four distinct morphological groups for the purpose of the current study, without currently being concerned

by the paleobiological/evolutionary status of each group. These morphogroups are briefly presented below.

Podocyrtis chalara Riedel and Sanfilippo

Pl. 1, fig. A, B

- 1970 Podocyrtis (Lampterium) chalara Riedel and Sanfilippo, p. 535, pl. 12, figs. 2, 3.
- 1971 *Podocyrtis (Lampterium) chalara* Riedel and Sanfilippo: Moore, p. 743, pl. 3, figs. 5,
 6.
- 1972 Lampterium chalara Riedel and Sanfilippo: Petrushevskaya and Kozlova, p. 543, pl.32, fig. 12.
- 1978 *Podocyrtis (Lampterium) chalara* (Riedel and Sanfilippo): Riedel and Sanfilippo, p.71, pl. 8, fig. 3, text-fig. 3.
- 2012 *Podocyrtis (Lampterium) chalara* Riedel and Sanfilippo: Kamikuri, p. 103, pl. 3, figs.
 2a, 2b.
- 2012 *Podocyrtis (Lampterium) chalara* Riedel and Sanfilippo: Moore and Kamikuri, p. 9, pl. P7, fig. 8.

<u>Distinguishing characters</u>: We include here forms displaying twelve or less vertically wellaligned, subangular abdominal pores of similar size per horizontal row, illustrating the classic morphology of *P. chalara*. Specimens of this morphogroup display less than thirteen pores on the circumference of the abdomen.

Podocyrtis sp. cf. P. chalara Riedel and Sanfilippo

Pl. 1, fig. C, D

- 1972 Lampterium sp. G: Petrushevskaya and Kozlova, pl. 32, fig. 10.
- 1972 Lampterium sp. aff. L. goetheana: Petrushevskaya and Kozlova, pl. 32, fig. 13.
- 2022 *Podocyrtis (Lampterium) chalara* Riedel and Sanfilippo: Meunier and Danelian, p. 21, pl. 2.4.

<u>Distinguishing characters</u>: This morphogroup includes specimens that have a similar outline and appearance as *P. chalara*, but they display vertically misaligned subangular abdominal pores of different size. Specimens of this morphogroup may display vertical rows of pores that are shifted to the right or left compared to the rows of pores above and below, giving a twisted appearance for the rows of pores developed on the abdomen, with result a honeycomb-like pore pattern. These shifts may be so extensive that the arrangement of pores on the abdomen becomes chaotic, preventing the possibility to trace any apparent abdominal rows or vertical alignment.

Podocyrtis sp. cf. P. goetheana (Haeckel)

Pl. 1, figs. E - L

2006 *Podocyrtis (Lampterium) chalara* Riedel and Sanfilippo: Funakawa et al., p. 29, pl.P9, figs. 11a, 11b.

<u>Distinguishing characters</u>: This morphogroup is mainly characterized by an increase in total size, but with a significant reduction in the number of abdominal pores compared to both variants of *P. chalara*. It differs from *P. goetheana* in that the bars of the second row of abdominal pores are thicker and not always elongated, nor parallel to each other, as the

formation of the honeycomb-like pattern of pores becomes more apparent. This morphogroup displays a high degree of morphological variability.

Podocyrtis goetheana (Haeckel)

Pl. 1, figs. M - O

- 1887 Cycladophora goetheana Haeckel, p. 1376, pl. 65, fig. 5.
- 1970 Podocyrtis (Lampterium) goetheana (Haeckel): Riedel and Sanfilippo, p. 535.
- 1971 Podocyrtis (Lampterium) goetheana (Haeckel): Moore, p. 743, pl. 3, figs. 7, 8.
- 1972 Lampterium sp. aff. L. goetheana Petrushevskaya and Kozlova, pl. 32, fig. 14.
- 2005 Podocyrtis (Lampterium) goetheana (Haeckel): Nigrini et al., p. 45, pl. P5, figs. 11,
 12.
- 2006 Podocyrtis (Lampterium) goetheana (Haeckel): Funakawa et al., p. 29, pl. P9, figs.12a, 12b.
- 2012 Podocyrtis (Lampterium) goetheana (Haeckel): Kamikuri, p. 103, pl. 3, fig. 1.
- 2012 *Podocyrtis (Lampterium) goetheana* (Haeckel): Moore and Kamikuri, p. 9, pl. P7, fig.
 9.

2022 Podocyrtis (Lampterium) goetheana (Haeckel): Meunier and Danelian, p. 21, pl. 2.5.

<u>Distinguishing characters</u>: This group includes only forms that display elongated straight bars formed at the level of the second horizontal row of pores on the abdomen. This feature is typical for *P. goetheana* as originally described, and, as mentioned above, the first occurrence of typical *P. goetheana* defines the base of the RP16 Zone.

Plate 3.1. Composite light micrographs of *Podocyrtis* radiolaria from ODP Site 1260, processed and scaled in ImageJ. (**A**) and (**B**) *Podocyrtis chalara*, samples: ODP 1260A-6R-5W, 63-65 cm and ODP 1260A-6R-5W, 20-

22 cm; (**C**) and (**D**) *Podocyrtis* sp. cf. *P. chalara*, samples: ODP 1260A-6R-4W, 68-70 cm and ODP 1260A-6R-5W, 20-22 cm; (**E**) to (**L**) *Podocyrtis* sp. cf. *P. goetheana*, samples: ODP 1260A-6R-4W, 68-70 cm; ODP 1260A-6R-5W, 15-17 cm and ODP 1260A-6R-5W, 87-89 cm; (**M**) to (**O**) *Podocyrtis goetheana*, sample: ODP 1260A-6R-1W, 58-60 cm.


3.3 Materials and methods

3.3.1 Sediment samples

The material analyzed in this study consists of radiolaria obtained from an expanded middle Eocene siliceous chalk sequence drilled at ODP Sites 1259 and 1260 (Leg 207, Demerara Rise), located in the equatorial region of the Atlantic Ocean, 380 km offshore Suriname (Erbacher et al., 2004; Danelian et al., 2005). The middle Eocene sequence is particularly thick at Sites 1259 and 1260 and contains siliceous microfossils (radiolarians, diatoms) of an excellent state of preservation (Danelian et al., 2007; Renaudie et al., 2010; Meunier and Danelian, 2022, 2023). The part of the limestone sequence from Site 1260 that is studied here is dated by orbito-chronology (Westerhold and Röhl, 2013). More specifically, our samples span the interval between 41.24 Ma and 39.84 Ma. Site 1259 is dated via bio- and magneto-stratigraphy and was sampled in the interval between ~39.05 and 37.70 Ma.

3.3.2 Slide preparation

A combined total of 15 samples from both sites were chosen and prepared for microscopic observation using techniques described by Sanfilippo et al. (1985). A small quantity (~2 cm³) of unprocessed sediment was collected from each sample and dried overnight at 50°C to eliminate any residual water. After being weighed, sediment samples were soaked for 2 hours in a 500 mL polypropylene beaker containing 30 mL of 30 % hydrochloric acid (HCl), to dissolve their carbonate content and concentrate siliceous microfossils. A few mL of HCl were added at the end to confirm the end of the reaction. The residues resulting from the acid treatment was then washed by adding ~200 mL of distilled water. After 2 hours of decantation, excess water was carefully removed using a pipette. Residues were subsequently soaked for 2 hours in 30 mL of 10 % hydrogen peroxide (H₂O₂) to remove organic matter, and subsequently

washed through a 63 μ m sieve using distilled water. The > 63 μ m fraction was then exposed to ultrasonic waves for 10 min, then passed again through the 63 μ m sieve, and finally left to dry overnight at 50°C. For each sample, ~2 g of dried residue was carefully spread on top of a slide covered with several drops of Norland Optical Adhesive 61, then topped with a coverslip and sealed by two minutes of exposure to UV light.

3.3.3 Microscopy and image processing

The resulting slides were analyzed with a Zeiss AXIO Images A2 microscope under transmitted light at ×10 and ×20 magnifications. All specimens recognized qualitatively as *P. chalara*, *Podocyrtis* sp. cf. *P. chalara*, *Podocyrtis* sp. cf. *P. goetheana* or *P. goetheana* were manually photographed using the mounted Axiocam ERc5s with Zen 3.5 (blue edition) software. For each specimen, a batch of 5-10 photographs were taken at different focal points to obtain a series of images, which were stacked afterwards using Helicon Focus 7.7.0 (HeliconSoft) to create a composite picture entirely in focus. The stacked images were subsequently retouched with Paint3D to facilitate their automated segmentation. This last procedure was performed with the ImageJ BioVoxxel plugin (Brocher, 2022) and the AutoRadio_Segmenter plugin developed for ImageJ / Fiji (Tetard et al., 2020).

3.3.4 Morphometric analyses

As we aimed to subject specimens to discriminant analysis based on morphometric measurements, it was essential that the documented morphological variables would adequately capture shape variations present in between the four morphogroups. We first designed a set of variables that would allow to compare with the variables used by Watanabe et al. (2022) on the specimens of the *Lampterium* lineage from the Pacific Ocean. A subset of seven of these

variables was retained and supplemented with six newly proposed variables to result in a set of thirteen morphological variables that document well the morphological variation between the four morphogroups (Figure 3.1). The seven morphological variables proposed by Watanabe et al. (2022) are:

- W1: Maximal width of the cephalis
- W2: Maximal width of the thorax
- W3: Maximal width of the abdomen
- H1: Maximal height of the cephalis
- H2: Maximal height of the thorax
- H3: Maximal height of the abdomen
- TL: Total length or height of the specimen

and our six additional variables are:

- LP2: Maximum length of the second abdominal pore along along the axis used to measure H3 (with the first pore being the one closest to the thorax-abdomen border)
- NPV: Number of abdominal pores aligned vertically along the axis used to measure H3 (on the front-facing side of the skeleton)
- NPH: Number of abdominal pores aligned horizontally behind the axis of W3 (on the front-facing side of the skeleton)
- R1: Maximum length of the second pore of the abdomen / H3
- R2: Number of abdominal pores aligned vertically behind the axis of H3 / H3
- R3: Number of abdominal pores aligned horizontally behind the axis of W3 / W3



Figure 3.1. Schematic representation of the eight skeletal variables measured for all specimens outlined in Table 3.1. Five additional variables related to pore counts and ratios are not illustrated. Abbreviations: W1: maximum width of the cephalis; W2: maximum width of the thorax; W3: maximum width of the abdomen; H1: maximum height of the cephalis without the apical horn; H2: maximum height of the thorax; H3: maximum height of the specimen without the apical horn; LP2: maximum length of the second abdominal pore.

Quantifications of these variables (Table 3.S1 in Supplementary materials) were directly performed on a dataset of 214 photographs/specimens from samples 1260A-6R-1W, 58-60 cm, 1260A-6R-4W, 68-70 cm, 1260A-6R-5W, 15-17 cm and 1260A-6R-5W, 87-89 cm, outlined in Table 3.1, using the image processing and analysis software ImageJ (Schneider et al., 2012). All of these data were tabulated and then imported in the statistical environment R (v. 4.1.3; R Core Team 2022) for subsequent LDA using the packages MASS (v. 7.3-60; Ripley et al., 2013)

and vegan (v. 2.6-4; Oksanen et al., 2013). This analysis is a constrained ordination procedure that uses a linear combination of coefficients to maximize the distance between *a priori* defined groups, while minimizing the distance within each group (Venables and Ripley, 2002). As morphometric variables were expressed in both metric units and as ratios, the LDA was performed on data that were transformed to have the mean at zero- and one-unit standard deviation (z transformation). After subjecting the whole dataset to LDA we performed cross-validation by 100 replicates of randomly assigning 80 % of the data to a training dataset and the remainder to a testing set to evaluate the classification success of the LDA.

Samples	P. chalara	Podocyrtis sp. cf. P. chalara	Podocyrtis sp. cf. P. goetheana	P. goetheana	Total number of measured specimens per sample
1260A 6R 1W 58-60cm	0	0	2	29	31
1260A 6R 4W 68-70cm	18	19	37	0	74
1260A 6R 5W 15-17cm	19	12	7	0	38
1260A 6R 5W 87-89 cm	34	9	28	0	71

Table 3.1. Number of specimens analyzed per morphogroup and per sampled core interval.

3.3.5 Artificial Neural Networks

Artificial Neural Networks or simply neural networks are machine learning algorithms designed to simulate the decision-making processes of the human brain by analyzing and exploiting patterns in data (Yang and Yang, 2014). Prior to analysis, the data given to a network is split into two parts, one for training and one for testing, usually in a 80 to 20 ratio. The first set of data is used to train the neural network, so that to enable it to learn recognizing features and patterns present in the data, whereas the second set of data is used to test the performance of the network to classify cases based on the previously trained capabilities of recognition.

A type of ANNs that has been specifically designed to analyze visual data are CNNs, commonly used for image recognition. They are designed to analyze visual data by considering the color values of each pixel and by identifying patterns within images (Hijazi et al., 2015). CNNs utilize a process known as convolution. Convolution can be described as a linear operation to decompose the input image by sliding small windows known as filters or kernels over the input image to construct layers that each obtain certain features. The convolutional layers in a CNN modify gradually the image parameters, such as weights or bias, to learn and recognize specific patterns or objects in the images. By adjusting these parameters through training, the network aims to correctly classify the output given a particular input. When a CNN has multiple layers, typically more than three, the procedure is referred to as deep learning, as each layer enables the recognition of more and more advanced features in an image. As mentioned, SNNs consider additionally the time factor, alike biological neurons, which use discrete spikes to compute and transmit information, instead of characterizing neurons by a single, static continuous-valued activation.

The hyperparameters (i.e. weights) of the neural network analyzed in this study were chosen randomly; more specifically a value between -1 and 1 was chosen for the VGG16 and SuperSpike-based networks, while a value network between 0 and 1 was chosen for the STDP-based network neural. We analyzed two sets of stacked and segmented images. The first set contained images with *a priori* assignments to the four morphogroups represented in section 2; for the second set assignments were altered based on the results of LDA. In each case, we performed ten runs per type of neural network used, with 20 epochs for each network, except

for the STDP-based network that was run with 100 epochs. An epoch simply means how many passes it goes through the training set and updates parameters based on each pass. The following neural networks were used to perform runs:

- Visual Geometry Group 16 (VGG16). This 16-layer deep CNN was used for its simplicity. For our analyses, we use transfer learning, meaning that the first 15 layers were already pre-trained based on a large-scale image dataset from ImageNet, and we only trained the last layer specifically using our data and PyTorch (Paszke et al., 2019). More information about VGG16 is provided by Simonyan and Zisserman (2015).
- STDP-based Spiking Neural Network (STDP-Network). This network contains convolutional and pooling layers that learn the features from the data using a Spike Timing Dependent Plasticity (STDP), a learning algorithm inspired by natural neurons. STDP adapts the synaptic connections between the neurons based on the timing of the spikes to transmit information (Masquelier and Thorpe, 2007) in an unsupervised way, i.e., without *a priori* group assignments. This SNN is then combined with a Support Vector Machine (SVM) for classification in the STDP-based network using the *a priori* group assignments (Cortes and Vapnik, 1995). To train the STDP-Network, we used the CSNN-simulator (Falez, 2019).
- SuperSpike-based Spiking Neural Network (SuperSpike-Network). This SNN is trained using a nonlinear voltage-based three-factor learning rule capable of training multilayer networks called the SuperSpike (Zenke and Ganguli, 2018), which is a supervised global learning rule similar to deep learning. We used the Norse simulator for our analyses (Pehle and Pedersen, 2021).

All the neural network simulations were conducted on the cluster "grouille" of the Grid'5000 test bed (Balouek et al., 2013) using two Nvidia A100-PCIE-40GB GPUs, an AMD EPYC 7452 32 core CPU (Zen 2, 2 CPUs/node), and 128GB of RAM. For each type of neural network, we averaged the assignment accuracies obtained over the ten replicate runs to gain robust insight into performance.

The training of the neural networks are expected to be better when a large, data-rich training set is used. Because our two analyzed datasets (Tables 3.2 and 3.3) are composed of 428 and 514 original images, respectively, we considered it necessary to augment the data available for ANN training. Therefore, we performed the following data augmentation procedures:

- Rotate the images by a randomized angle between -15 and 15 degrees and keep all copies.
- Randomly choose images that would be flipped from left to right and keep both copies.
- Rescale random images with values between 1 and 1.3 (with 1 being the default scale value).

The total number of images in each dataset was enhanced to >1000 images via data augmentation.

 Table 3.2. List of images included (prior to augmentation) in each of the classes for the four- class dataset used

 for analysis withVGG16with.

Samples	P. chalara	Podocyrtis sp. cf. P. chalara	Podocyrtis sp. cf. P. goetheana	P. goetheana	Total number of analyzed specimens per sample
1260A 6R 4W 68-70cm	0	19	34	3	56
1260A 6R 4W 119- 121cm	0	13	2	1	16
1260A 6R 5W 15-17cm	0	12	7	0	19

TOTAL	174	148	87	18	427
1260A 8R 3W 65-67cm	0	1	0	0	1
1260A 7R 2W 19-21 cm	0	13	0	0	13
1260A 7R 1W 121-123 cm	31	11	0	0	42
1260A 7R 1W 69-71 cm	33	11	0	0	44
1260A 7R 1W 22-24 cm	24	15	0	0	39
1260A 6R 6W 57-59 cm	19	14	2	0	35
1260A 6R 6W 20-22 cm	12	18	9	0	39
1260A 6R 5W 87-89 cm	33	12	18	14	77
1260A 6R 5W 63-65cm	22	9	15	0	46

 Table 3.3. List of images included (prior to augmentation) in each of the classes for the three class three dataset

 used for analysis with VGG16, STDP- based SNN and SuperSpike-based SNN.

Samples	P. chalara + Podocyrtis sp. cf. P. chalara	Podocyrtis sp. cf. P. goetheana	P. goetheana	Total number of analyzed specimens per sample
1259A 17R 1W 54-56cm	0	0	1	1
1259A 18R 1W 53-55cm	0	0	41	41
1259A 18R 2W 53-55cm	0	0	36	36
1260A 6R 1W 58-60cm	0	1	26	27
1260A 6R 4W 68-70cm	19	29	0	48
1260A 6R 4W 119- 121cm	13	2	1	16
1260A 6R 5W 15-17cm	12	5	0	17
1260A 6R 5W 63-65cm	31	15	0	46
1260A 6R 5W 87-89cm	43	28	0	71
1260A 6R 6W 20-22cm	29	8	0	37
1260A 6R 6W 57-59cm	34	2	0	36
1260A 7R 1W 22-24cm	39	0	0	39
1260A 7R 1W 69-71cm	44	0	0	44
1260A 7R 1W 121- 123cm	42	0	0	42
1260A 7R 2W 19-21cm	13	0	0	13
TOTAL	319	90	105	514

3.4 Results

3.4.1 Morphometrics and linear discriminant analysis

The LDA performed on the matrix of our 13 morphometric variables, i.e., measurements, pore counts and ratios, represented >99 % of the variation on the first two axes and clustered *Podocyrtis goetheana* and *Podocyrtis* sp. cf. *P. goetheana* successfully (Fig. 3.2). Comparatively, specimens belonging to *P. chalara* were regularly confused with *Podocyrtis* sp. cf. *P. chalara* and *vice versa*, resulting overall in 73.5 ± 6.1 % (mean \pm sd) correct assignments (Table 3.4). These two latter morphogroups overlap completely on the LDA plot, whereas *P. Podocyrtis* sp. cf. *P. goetheana* and *P. goetheana* are relatively well-separated from *P. chalara* and *Podocyrtis* sp. cf. *P. chalara*, although they share limited overlap with each other (Figure 3.2). The first axis of the LDA mainly represents size variations of the second row of abdominal pores. The morphological variables that contribute the most to discriminate the morphogroups on the first LDA axis correspond to the maximum length of the second abdominal pore (LP2), and the maximum length of the second pore of the abdomen/maximum height of the abdomen (R1).



Figure 3.2. Scatter plot of the two first axes of the Linear Discriminant Analysis (LDA) conducted on the 13 variables that constituted the morphometric data using four *a priori* identified morphogroups (indicated in the legend; ellipses correspond to the 95 % confidence intervals for each morphogroup). Abbreviations: W1: maximum width of the cephalis without the apical horn; W2: maximum width of the thorax; W3: maximum width of the abdomen; H1: maximum height of the cephalis; H2: maximum height of the thorax; H3: maximum height of the abdomen; TL: total length or height of the specimen without the apical horn; LP2: maximum length of the second abdominal pore; NPV: number of abdominal pores aligned vertically; NPH: number of abdominal pores aligned horizontally; R1: maximum length of the second pore of the abdomen/H3; R2: number of abdominal pores aligned horizontally/W3.

Upon considering three morphogroups achieved by lumping *P. chalara* and *Podocyrtis* sp. cf. *P. chalara*, the LDA ordination is highly similar to that obtained in the four-group analysis, with all variation represented on the first two axes (~92.87 % on axis 1; **Figure 3.S1**). Classification results improved substantially, with 94.5 \pm 3.1 % of correctly identified

specimens in the three-group LDA (Table 3.S2). Misclassification mainly occurred between *P. goetheana* and *Podocyrtis* sp. cf. *P. goetheana*, occasionally also between *P. chalara* and *Podocyrtis* sp. cf. *P. goetheana*, but never between *P. chalara* and *P. goetheana*.

 Table 3.4. Average confusion matrix of the *Podocyrtis* morphotypes based on linear discriminant analysis (LDA)

 performed on the matrix of our 13 morphometric variables (.,i.e., measurements, pore counts and ratios).

LDA results 4-group								
scenario								
	P_sp_cf_P_cha_a		P_cha_and_	P_goe_and_P		P_goe_a	total_	
	nd_P_sp_cf_P_go	P_cha_and_P	P_sp_cf_P_g	_sp_cf_P_ch	P_goe_and_P	nd_P_ch	correc	
	e	_sp_cf_P_cha	oe	a	_sp_cf_P_goe	a	t	
res_								
mea							0.7351	
n	0.011162791	0.204651163	0.001627907	0	0.04744186	0	16279	
res_							0.0612	
sd	0.013819733	0.059589621	0.005963544	0	0.024717954	0	97375	

3.4.2 Artificial neural networks

3.4.2.1 Classification using the four morphogroups

First, we trained neural networks on the dataset of 428 images attributed to the four morphogroups (Table 3.2 4_classes dataset) using a CNN with a VGG16 architecture. Prior to the training and testing phases, images were manually grouped into four distinct classes (one class per morphogroup). The analysis was run ten times with the pre-processing parameters outlined in Table 3.5, which resulted to an identification accuracy of 54.4 ± 1.7 % (Figure 3.3).These results indicate that, although the network was able to partially identify the differences between the general morphologies of *P. chalara* + *Podocyrtis* sp. cf. *P. chalara versus* those of *P. goetheana* + *Podocyrtis* sp. cf. *P. goetheana*, it was not able to accurately distinguish all four morphogroups from each other.

Table 3.5. Summary of the tests performed on the 4_classes and 3_classes datasets using VGG16, STDP-based SNN and SuperSpike-based SNN. Accuracy values are averaged from the results of the respective 10 test runs and indicated as the mean \pm one standard deviation.

Dataset	Architecture	Test runs	Number of Epochs	Data pre-processing	Simulation time h:m:s	Accuracy (%) ± std mean ± std
4_classes	VGG16	10	20	Image resize to (224, 224) px	0:31:46	54.40 ± 1.74
	VGG16	10	20	Image resize to (224, 224) px	0:17:37	92.60 ± 0.77
3_classes	STDP-Network	10	100	Image resize to (128, 128) px and On-Off filter	0:05:4343	90.40 ± 0.4949
	SuperSpike- Network	10	20	Image resize to (128, 128) px	0:23:26	84.42 ± 1.36



Figure 3.3. Example of a confusion matrix obtained from random single run – Analysis of 4_classes dataset with VGG16. Correct assignments were reached in 54.40 % of the cases. The color scale indicates the number of specimens.

3.4.2.2 Classification using the three morphogroups supported by LDA

Subsequently, we trained neural networks using a dataset with *a priori* assignment to the three morphogroups that were recognized by LDA, i.e., *P. chalara* + *Podocyrtis* sp. cf. *P. chalara Podocyrtis* sp. cf. *P. goetheana* and *P. goetheana*. For these analyses we compared the performance of a VGG16 CNN, a STDP-based SNN and a SuperSpike-based SNN using the 3_classes dataset (Table 3.3). The resulting assignment results and average network accuracies

(Figure 3.4) indicate that, under our specified conditions, all neural networks are able to accurately assign specimens to their correct class when three predefined classes are used in combination with a large dataset of images. However, we observed substantial differences in the speed to conclude analyses; thus the STDP-based SNN was the fastest (around 6 minutes) due to the use of the local learning rule for training and one spike per image per neuron principale. Moreover, VGG16 came second with a time of around 17 minutes due to the size of the network and the use of transfer learning and training of only the last layer. Last came the SuperSpike-based SNN with approximately 23 minutes per run because all layers of the network were trained from scratch during the training phase at the start of each run using a global learning rule.



Figure 3.4. Examples of confusion matrices obtained from random single runs – Analyses of 3_classes dataset with VGG16, STDP-based SNN and SuperSpike-based SNN. Correct assignments were reached in 92.60 %, 90.40 % and 84.42 % of the cases in VGG16, STDP-based SNN and SuperSpike-based SNN, respectively. The color scale indicates the number of specimens.

3.5 Discussion

In this study, we examined the morphological variability in Eocene *Podocyrtis* belonging to the anagenetic sequence that starts with *P. chalara* and ends with *P. goetheana*. Specifically, we examined and compared the performance of two analytical approaches in assigning individuals to *a priori* defined morphogroups that were constructed from qualitative observations, i.e., *P. chalara*, *Podocyrtis* sp. cf. *P. chalara*, *Podocyrtis* cf. *P. goetheana* and *P. goetheana*. The first approach involved LDA based on morphometric data, whereas the second was a neural network approach based on automatic image recognition. Both methods gave very similar results, which indicates that both morphometrics and image analysis evaluated shape differences in a highly similar way, suggesting that the results obtained with these two different methods are robust.

Comparing LDA and neural network approaches, assignment performances were comparatively low when four morphogroups were considered in the morphological transition from *P. chalara* to *P. goetheana*. The scatterplot of the LDA (Figure 3.2), indicated a large morphospace overlap between *P. chalara* and *Podocyrtis* sp. cf. *P. chalara*, revealing that the qualitatively observed 'differences' were either not sampled in our datasets or that these differences are part of a larger spectrum of morphological variation and not informative to distinguish morphogroups. Given the similarity of LDA and neural networks, based on different datasets, the second hypothesis is more likely; it appears that *P. chalara* and *Podocyrtis* sp. cf. *P. chalara*. Finally, *Podocyrtis* sp. cf. *P. goetheana* also overlaps with *P. chalara* and *Podocyrtis* sp. cf. *P. chalara* in morphospace occupation, however this overlap is small. The analyses conducted with the machine learning

approach on four morphogroups confirm the results of LDA, as VGG16 had significant difficulties in differentiating specimens of *P. chalara* from *Podocyrtis* sp. cf. *P. chalara* and *vice-versa* using stacked and segmented images, when all four morphogroup classes were predefined. This was also the case for some specimens of *Podocyrtis* sp. cf. *P. goetheana* and *P. goetheana*, which all resulted in an inferior performance of VGG16 compared to that of the LDA for the scenario with four morphogroups,, i.e., 54.4 ± 1.7 % versus 73.5 ± 6.1 %, respectively.

When *P. chalara* and *Podocyrtis* sp. cf. *P. chalara* are lumped in the same morphogroup, resulting in a three-group configuration, assignment probabilities improved strongly both for LDA and neural networks. On average 94.5 ± 3.1 % of assignments were correct with LDA, whereas 92.6 ± 0.8 %, 90.4 ± 0.5 % and 84.4 ± 1.4 % of the assignments were correct for VGG16, STDP-based SNN and SuperSpike-based SNN, respectively. Neural networks, mainly our CNN can accurately and quickly, assign specimens to the three pre-defined classes using a large dataset. Whereas the CNN performed highly similarly to LDA, both SNNs we used here performed less well, as had been documented for other tasks before (Tavanaej et al., 2019). Runs with the STDP-based SNN ran to completion fastest and given that the accuracy was only slightly reduced compared to LDA and VGG16, this approach may be preferred for datasets that require a very long runtime with similar CNNs. The accuracy of the SuperSpike-based SNN was significantly reduced compared to all the other classification methods that we used. Further work is required to determine the cause of this underperformance, but the lower accuracy for the SuperSpikeSuperSpike-based SNN is possibly due to the network size, as we used only using eight layers in our case compared to 16 layers in VGG-16.

As for the performance evaluation of the neural network approach, it is noteworthy that images were obtained manually for the purpose of our study, but advances in image technology now allow that much of the image acquisition and preparation procedures (photographing, stacking and segmentation) to be automated by the use of automatic microscopes and modification of the AutoRadio_Segmenter plugin's code (Marchant et al., 2020; Tetard et al., 2020). Using these automated procedures would facilitate the scalability of the entire analysis with larger datasets. If such automated procedures were to be used, constructing image datasets may potentially become more time-efficient than the various procedures that are required to develop a morphometric dataset. Another advantage of using neural network is their quick run time (Table 3.5.), although both accuracy and run time would increase upon using larger image datasets. Furthermore, standardizing the rotation and orientation of specimens is essential in morphometric studies; however, this requirement can be relaxed for neural networks, as in our case re-rotating and re-orienting were used in the data augmentation process to enlarge input datasets.

In our study, we evaluated assignment accuracy of LDA and machine learning with neural networks based on *a priori* group assignments; however, in the future, it would be useful to examine morphological variation without considering such assignments, e.g., by using other ordination techniques and/or by using unsupervised machine learning techniques. These methods could help in attempts to evaluate whether the three retained morphogroups represent natural entities, although based on fossil evidence only such assessments are very difficult. We refrain from such analyses here, as we believe these would be best conducted with a larger set of material that ideally would cover the total geographic range and the total stratigraphic interval covered by *P. chalara*, *Podocyrtis* sp. cf. *P. goetheana* and *P. goetheana*.

Future work could also be focused on trying to evaluate the capacity of VGG16 to accurately differentiate between the P. chalara and Podocyrtis sp. cf. P. chalara specimens, with an altered set of images. Indeed, the weighted activation heatmap (Grad-CAM) (Servaraju et al., 2017) generated from the runs with the 3-class dataset showed that the analytical focus of the network was centered around the thoracic and abdominal walls of the specimens (Figure 3.S2 in Supplementary material). One could try to develop a dataset including several unstacked and non-segmented images per specimen, each one with a focus on specific morphological features. This would allow some images to contain as much detail as possible on abdominal features, whereas others would focus on other regions (e.g. the thorax) and include blurred features of the outer walls and backside of the abdomen. We hypothesize that these manipulations could force the network to focus its recognition capabilities on a larger set of morphological features. If this hypothesis would be correct, it would also facilitate the data acquisition process by eliminating the need to stack and segment the images themselves. Alternatively, it is possible that providing more fragmented information to neural networks would hamper an efficient learning process, with negative consequences on the accuracy of the following predictions, somewhat similar to what we observed in the VGG16 evaluation based on the 4-class dataset.

3.6 Conclusion

The aim of our work was to study the morphological variability in the anagenetic sequence of *P. chalara* to *P. goetheana* and to evaluate the performance of recognizing and classifying four *a priori* identified morphogroups with various machine learning algorithms based on neural networks that use image data as direct input in comparison to linear discriminant analysis using morphometric data. Our results demonstrate that LDA and neural networks provide very similar outcomes, indicating robust performances. With both approaches we encountered difficulties

distinguishing *P. chalara* and *Podocyrtis* sp. cf. *P. chalara*, suggesting that the qualitative basis on which these morphogroups were recognized is to be revised. For both approaches, assignment probabilities drastically increased for the scenario where three morphogroups were recognized, lumping P. chalara and Podocyrtis sp. cf. P. chalara, and thus retaining P. chalara + Podocyrtis sp. cf. P. chalara, Podocyrtis sp. cf. P. goetheana and P. goetheana. Further studies with more comprehensive sampling are required to document the likelihood of Podocyrtis sp. cf. P. goetheana representing a separate natural entity, which additionally has implications for the position of bioevent RP16. However, our results indicate that the morphometric data that we used for LDA samples the morphological variability in the anagenetic sequence of *P. chalara* to *P. goetheana* in a comprehensive way. Secondly, neural network approaches were able to correctly assign most specimens, and therewith to accurately distinguish the three morphogroups directly from specimen images. These results indicate that VGG16, STDP-based SNNs, and even SuperSpike-based SNNs are capable of recognizing morphological variation in images and thus of reliably distinguishing radiolarian morphogroups, which could fascilitate identification and help with reaching more objective taxonomic decisions. Furthermore, neural network approaches can be combined with automated image acquisition and preparation procedures (photographing, stacking and segmentation) that enable the creation of much larger image databases in a time-efficient manner. Analyses based on neural network architecture could thus take a fraction of the time that would be required for a trained taxonomist/(paleo)biologist to create and analyze quantitative morphometric datasets. In conclusion, neural network approaches based on images of (paleo)biological specimens may provide promising opportunities to guide more objective taxonomic decisions.

Data availability

Microscopic slides are prepared and stored at UMR 8198 – Evo-Eco-Paleo of the University of Lille, France. The datasets (https://doi.org/10.57745/8KBOFP, Pinto et al. 2023) have been archived in the repository of the University of Lille at Recherche Data Gouv. The codes are available at

https://archive.softwareheritage.org/browse/directory/cc7d8ef1505299a208adcde597a98d90b Oca47d6/https://archive.softwareheritage.org/browse/directory/cc7d8ef1505299a208adcde597 a98d90b0ca47d6/ (Elbez, 2023).

Acknowledgments

This study was partly funded by the French government through the program "Investissements d'avenir" (I-ULNE SITE/ANR-16-IDEX-0004 ULNE) managed by the National Research Agency. It also received funding from the European Union's Horizon 2020 research and innovation program under the Marie Sklodowska-Curie grant agreement no. 847568. It was also supported by UMR 8198 Evo-Eco-Paléo and IRCICA (CNRS and Univ. Lille USR-3380).

A special thanks to Mazdak Fatahi for very helpful discussions regarding machine learning techniques and to Sylvie Regnier for help and advice with the sample preparation.

PART 2: ARTIFICIAL INTELLIGENCE APPLIED TO ASSEMBLAGES

<u>CHAPTER 4</u>: Subchapter 1 - Convolutional neural network application on a new middle Eocene radiolarian dataset

Veronica Carlsson^{1,2}*, Taniel Danelian¹, Martin Tetard³, Mathias Meunier¹, Pierre Boulet², Philippe Devienne², Sandra Ventalon⁴

¹Univ. Lille, CNRS, UMR 8198, Evo-Eco-Paleo, F-59000 Lille, France.

²Univ. Lille, CNRS, CRIStAL – Centre de Recherche en Informatique Signal et Automatique

de Lille, UMR 9189, F-59000 Lille, France.

³GNS Science, NZ-5040, Lower Hutt, New Zealand.

⁴Univ. Lille, CNRS, Univ. Littoral Cote d'Opale, UMR 8187, LOG, F-59000 Lille, France.

* Corresponding author. E-mail address: veronica.carlsson@univ-lille.fr

Published in Marine Micropalaeontology, 183, 2023 https://doi.org/10.1016/j.marmicro.2023.102268

Version presented here : Accepted (No Open Access)

Author contributions

Veronica Carlsson is the main author and writer of this paper. She formulated the research methodology of this study in discussion with the other authors. In addition, she classified, analyzed, collected, and processed all images, established the taxonomic catalogue, and built up the image datasets.

Taniel Danelian contributed to the formulation of the research methodology and with supervision and guidance of this study. He also contributed to the revision and discussion of the taxonomy and the writing and revision of the manuscript.

Martin Tetard contributed to the formulation of the research methodology and revision of the datasets, as well as to running the training data. He also contributed to the revision of the paper.

Mathias Meunier contributed to the taxonomic revision of the datasets and the writing and reviewing of the taxonomic catalogue.

Pierre Boulet contributed to the supervision of this study and the revision of the manuscript.

Philippe Devienne contributed to the supervision and revision of the manuscript.

Abstract

A new radiolarian image database was used to train a Convolutional Neural Network (CNN) for automatic image classification. The focus was on 39 commonly occurring nassellarian species, which are important for biostratigraphy.

The database consisted of tropical radiolarian assemblages from 129 middle Eocene samples retrieved from ODP Holes 1258A, 1259A, and 1260A (Demerara Rise). A total of 116

taxonomic classes were established, with 96 classes used for training a ResNet50 CNN. To represent the diverse radiolarian assemblage, some classes were formed by grouping forms based on external morphological criteria. This approach resulted in an 86.6% training accuracy.

A test set of 800 images from new samples obtained from Hole 1260A was used to validate the CNN, achieving a 75.69% accuracy. The focus then shifted to 39 well-known nassellarian species, using a total of 15 932 images from the new samples. The goal was to determine if the targeted species were correctly classified and explore potential real-world applications of the trained CNN.

Different prediction threshold values were experimented with. In most cases, a lower threshold value was preferred to ensure that all species were captured in the correct groups, even if it resulted in lower accuracies within the classes.

Keywords: middle Eocene, radiolaria, convolutional neural network, image database, automated identification, image recognition

4.1 Introduction

Polycystine radiolaria are microscopic unicellular protists living currently in all modern oceans; they are characterized by an aesthetically pleasing siliceous skeleton known in the fossil record since the Cambrian (Pouille et al., 2011; Aitchison et al., 2017). Their fossil record is thus of much interest for evolutionary studies (i.e., Danelian and Johnson, 2001; Renaudie and Lazarus, 2013; Tetard et al., 2017; Trubovitz et al., 2020). They are usually the only fossils capable of dating siliceous sedimentary sequences (i.e., Danelian et al., 2012; Vrielynck et al., 2003) and are commonly used in paleoceanography (Matsuzaki et al., 2018; Itaki et al., 2020). Due to the

small number of experts, radiolarian taxonomy is less well-elaborated than the one of other microfossil groups, such as foraminifera and nannofossils.

Today, most of the studies involving identification and counting of microfossils, such as radiolarians, are conducted manually and require substantial taxonomic expertise. This process is known to be time-consuming, particularly when microfossils are used for paleoceanography. Moreover, consistency in species classification may be difficult to achieve between different taxonomic experts. Therefore, artificial intelligence (AI) has been introduced to this field to simplify or automate the work done by micropaleontologists, as for example through automatic image recognition and counting. Several applications of CNNs for automatic image recognition were developed over the last 20 years, since the introduction of SYRACO by Dollfus and Beaufort (1999). Nowadays, CNNs are used on various microfossil groups, such as foraminifera (ex. Mitra et al., 2019, Hsiang et al., 2019; Marchant et al., 2020), coccoliths (ex. Dollfus and Beaufort, 1999; Beaufort and Dollfus, 2004), pollen (ex. Bourel et al., 2020), or even radiolarians (Itaki et al., 2020; Renaudie et al., 2018; Tetard et al., 2020).

Sediments recovered from the Demerara Rise (tropical Atlantic Ocean) during the Leg 207 are rich in middle Eocene radiolarians, preserved in a continuous and expanded carbonate sequence. The encountered radiolarian diversity is immense and based on our estimates it accounts for ca. 500 species, many of which are not described yet. Indeed, although Eocene radiolarians have been studied for about 150 years (since Ehrenberg, 1874) and more in depth for the last 50 years (Riedel and Sanfilippo, 1970, 1978), past research was mainly focused on their biostratigraphic applications (see Meunier and Danelian, 2022 and references therein).

Taking advantage of today's technological achievements, our objective was to design a reliable approach to automatically classify middle Eocene tropical radiolarians from Demerara Rise. The main question addressed in this study is whether a CNN can accurately classify 39 commonly observed nassellarian radiolarian species, most of which have an established biostratigraphic significance. We thus trained a CNN on a newly established image database of middle Eocene radiolarians with a focus on some common nassellarian species. To do this we classified every single object appearing on prepared radiolarian slides. We were inspired by a recent similar study, conducted by Tetard et al. (2020), who studied middle Miocene to Quaternary radiolarians from the West Pacific Warm Pool. We also included additional *Podocyrtis* species prepared for an earlier dataset (Carlsson, 2022). Finally, a new set of a small number of samples was imaged to test the consistency of our trained CNN, which was double checked with manual identifications made by a taxonomist.

4.2 Materials and methods

4.2.1 Core setting and sample preparation

The middle Eocene samples used in this study were collected during ODP Leg 207 from Demerara Rise, situated off the coast of Suriname (Erbacher et al., 2004, Danelian et al. 2005). This dataset includes samples from several cores recovered from sites 1258, 1259 and 1260. It is noteworthy that the middle Eocene sequence at site 1260 is thick and contains radiolarians of an excellent state of preservation (Danelian et al., 2007). The full sample list used in this study may be found in Supplements 4.1, Table S1. The sediment samples consist primarily of nannofossil and foraminifera chalk, but also contain abundant and well-preserved siliceous microfossils, composed essentially of radiolarians (Danelian et al. 2007, Meunier and Danelian, 2022), as well as diatoms (Danelian et al. 2007, Renaudie et al. 2010).

Sediment samples from ODP holes 1258A, 1259A, and 1260A were first processed to disaggregate organic matter and dissolve carbonates and were then sieved through a 45 μ m mesh to remove smaller particles. Thereafter, the samples were prepared using a recent random settling protocol described in Tetard et al. (2020). About 1/3 to ¹/4 of a microspoon spatula was used from the dried residues, corresponding to approximately 0.2-1.0 g for each sample. Samples were uniformly settled onto 12 mm x 12 mm cover slides using a 3D-printed decanter, as in Tetard et al. (2020), preventing contact between radiolarians and other remaining objects on the slide. The cover slides were allowed to dry overnight before being mounted with NOA81 optical glue. A total of eight different cover slips, all from the same samples, were placed onto one 76 mm x 26 mm glass slide, forming one sample.

Finally, a new set of four samples (see Supplements 4.1, Table S1) was prepared for manually testing the actual accuracy of the trained CNN. An improved cleaning technique was used, which kept only siliceous particles that are larger than $63 \mu m$, and completely removed all clay, calcite and smaller or broken radiolarians (Sanfilippo et al., 1985, Tetard et al., 2020). In fact, most radiolarians are larger, so using a $63 \mu m$ sieve will probably just remove smaller broken pieces or radiolarians rather than the radiolarians themselves. At first, about 2-3 cm of raw sediment sample was placed in a plastic beaker. Thereafter 30 ml of 30 % hydrochloric acid was added and left for two hours and until there was no more reaction. Furthermore, 200 ml of distilled water was added to the beaker, which was stirred gently and left to settle for two additional hours. The supernatant was removed and exchanged with 30 ml of 10 % hydrogen peroxide and was left to rest for another two hours. The residue was thereafter washed with a $63 \mu m$ sieve into a 100 ml beaker. To remove remaining clay particles, the samples were

processed in ultrasonic bath for ten minutes; they were later sieved again with a 63 μ m mesh and collected afterwards in a filter. Next, the residues were dried in an oven at 50 °C and transferred into a glass vial.

4.2.2 Image collection and processing

All samples were photographed using a Nikon Eclipse Ni automatic microscope equipped with a Nikon DS-Ri2 microscope camera and Nikon NIS Element software, using a 20 x objective, allowing a 200 x magnification and 0.36 μ m/pixel resolution. The lowest and highest focal points on the Z-axis were manually determined for each sample since the depth varied depending on the thickness of the glass, coverslip, optical glue, and individual radiolarian specimens. For each field of view (FOV), an image was taken at different focal depths, every 10 μ m, depending on the minimum and maximum focal points determined for each sample. The microscope then automatically stacked all images, taken at different focal points for each FOV, in order to create a composite image. The microscope was programmed to capture small images in 20 x 20 (400) FOV, covering about 10 x 10 mm out of the 12 x 12 mm available surface for each individual sample. The individual FOV images were subsequently merged into one large mosaic image (see Fig. 4.1), which has as a result to lose less images of specimens located on the edges of FOVs.



Figure 4. 1) For each sample, 20×20 images are automatically photographed in a convolving way and merged into one large "mosaic composed" image, which enables preserving more complete specimens which are not cut in half. 2) The mosaic composed image thereafter went through grayscale conversion. 3) Segmentation of each unique particle into vignettes and image conversion into 8-bit black and white with black background.

The composite mosaic pattern images received from the automatic microscope were first converted from RGB to 8-bit grayscale using Adobe Photoshop to decrease its size. The images were then segmented into ImageJ using the ImageJ BioVoxxel plugin (Brocher, 2022), and a modified script of the Autoradio_Segmenter plugin (Tetard et al., 2020), which enabled each individual particle to become its own individual image, or vignette. For more details, the reader is referred to Tetard et al. (2020).

Finally, we also included additional *Podocyrtis* species images prepared for an earlier dataset (Carlsson, 2022).

4.2.3 Taxa selection and dataset

For this study we decided to focus on 39 species (Plates 4.1 and 4.2), which are the most common in the Middle Eocene interval of Demerara Rise and most of which are used in biostratigraphy (Riedel and Sanfilippo, 1970, 1978; Sanfilippo and Nigrini, 1998; Meunier and Danelian 2022).

Plate 4.1. Nassellarian radiolarian species commonly occurring in Middle Eocene sediments of Demerara Rise; species names are followed by the ODP site and hole, core, section and sampled level from which it comes from. A) Dendrospyris stylophora (Ehrenberg, 1874) from 1259A-20R-4 W, 53-55 cm, B) Elaphospyris didiceros (Ehrenberg, 1874) group from 1258A-2R-4 W, 55-57 cm, C) Liriospyris clathrata (Ehrenberg, 1874) group from 1259A-20R-4 W, 53-55 cm, D) Rhabdolithis pipa Ehrenberg 1854 from 1260A-15R-1 W, 55-57 cm, E) Lophophaena radians (Ehrenberg, 1874) group from 1259A-16R-2 W, 55-57 cm, F) Dictyoprora mongolfieri (Ehrenberg, 1854) from 1260A-6R-2 W, 55-57 cm, G) Dictyoprora amphora (Haeckel, 1887) 1258A-2R-4 W, group from 55-57 cm, H) Rhopalosyringium? auriculaleporis (Clark and Campbell, 1942) from 1260A-14R-6 W, 55-57 cm, I) Rhopalosyringium? biauritum (Ehrenberg, 1874) from 1260A-12R-3 W, 55-57 cm, J) Dictyomitra parva (Kim, 1992) from 1258A-2R-4 W, 55–57 cm, K) Lithochytris vespertilio Ehrenberg, 1874 from 1260A-10R-5 W, 55–57 cm, L) Sethochytris triconiscus (Haeckel 1887) from 1259A-20R-4 W, 53-55 cm, M) Lychnocanium babylonis (Clark and Campbell 1942) group from 1258A-3R-3 W, 56-58 cm, N) Lychnocanoma bajunensis (Renz, 1984) from 1258A-2R-4 W, 55-57 cm, O) Stichopterygium microporum (Ehrenberg, 1874) from 1260A-8R-4 W, 54-56 cm, P) Carpocanopsis ornata (Ehrenberg, 1874) group from 1260A-6R-5 W, 55–57 cm, Q) Cycladophora spatiosa (Ehrenberg, 1874) group from 1259A-20R-1 W, 55– 57 cm, R) Anthocyrtis mespilus (Ehrenberg, 1847) group from 1259A-16R-2 W, 55-57 cm, S) Zealithapium mitra (Ehrenberg, 1874) from 1258A-2R-3 W, 55-57 cm.



Plate 4.2. Nassellarian radiolarian species commonly occurring in Middle Eocene sediments of Demerara Rise; species names are followed by the ODP site and hole, core, section and sampled level from which it comes from. A) Lophocyrtis alauda (Ehrenberg, 1874) from 1260A-15R-1 W, 55-57 cm, B) Aphetocyrtis zamenhofi Meunier and Danelian, 2023 from 1259A-26R-5 W, 54-56 cm, C) Theocyrtis scolopax (Ehrenberg, 1874) from 1260A-15R-3 W, 55-57 cm, D) Phormocyrtis embolum (Ehrenberg, 1874) from 1258A-2R-4 W, 55-57 cm, E) Phormocyrtis lazari Meunier and Danelian, 2023 from 1260A-8R-6 W, 54-56 cm, F) Podocyrtis (Lampterium) goetheana (Haeckel, 1887) from 1259A-18R-1 W, 53-55 cm, G) Podocyrtis (Lampterium) chalara Riedel and Sanfilippo, 1970 from 1260A-6R-CC, 63-177 cm, H) Podocyrtis (Lampterium) mitra Ehrenberg group, 1854 from 1260A-9R-1 W, 55-57 cm, I) Podocyrtis (Lampterium) sinuosa Ehrenberg, 1874 from 1259A-15R-1 W, 55–57 cm, J) Podocyrtis papalis Ehrenberg, 1847 from 1258A-2R-4 W, 55–57 cm, K) Podocyrtis (Podocyrtoges) ampla Ehrenberg, 1874 from 1260A-10R-5 W, 55-57 cm, L) Podocyrtis (Podocyrtoges) physis Sanfilippo and Riedel, 1973 from 1259A-16R-1 W, 55-57 cm, M) Podocyrtis (Podocyrtoges) diamesa Riedel and Sanfilippo, 1970 from 1259A-26R-3 W, N) Podocyrtis (Lampterium) puellasinensis Ehrenberg, 1874 from 1259A-20R-4 W, 53-55 cm, O) Calocyclas hispida (Ehrenberg, 1874) from 1260A-6R-4 W, 55-57 cm, P) Thyrsocyrtis (Thyrsocyrtis) rhizodon Ehrenberg, 1874 from 1260A-6R-CC, 63-177 cm, Q) Thyrsocyrtis (Pentalocorys) triacantha (Ehrenberg, 1874) from 1260A-8R-6 W, 1259A-25R-2 W, 54–56 cm, R) Eusyringium lagena (Ehrenberg, 1874) from 55–57 cm, S) Eusyringium fistuligerum (Ehrenberg, 1874) group from 1259 to 18R-1 W, 53-55 cm, T) Rhopalocanium ornatum (Ehrenberg, 1874) from 1259A-22R-1 W, 55-57 cm.



Synonymy lists are given in the supplementary catalogue (Supplements 4.3), thus allowing the reader to understand the species concept followed in this study. Taxonomic information for all the other radiolarian classes used in the analysis is also presented in there; most of the other radiolarians were grouped in supraspecific taxa, with taxonomic information and some typical forms given in the catalogue (Supplements 4.3).

The ParticleTrieur software version 2.4.10, developed by Marchant et al. (2020), was used to label our dataset. It includes a built-in k-nearest neighbor (KNN) algorithm, which is a machine learning algorithm that can be used for supervising the classification tasks. It identifies the k-nearest training data points or neighbors for a new data point and predicts a label for the new data based on already labeled data. In the context of image classification, the data points correspond to the pixels in the image. Therefore, ParticleTrieur can recognize patterns in the images for individual classes and suggest classification for new or unclassified images, after some classes have already been added in a semi-supervised way but have of course been validated by a human expert.

We managed to build a dataset consisting of 12,217 images out of a total of ca. 50 000 images, distributed in 116 classes, including the 39 important key-classes of well-known nassellarian species (Plates 4.1 and 4.2). Some of the classes consisted of as few as one specimen per class, while others contained up to nearly a thousand images (i.e. the largest class). Classes with fewer than ten specimens were excluded from the CNN training, resulting in only 96 classes to be trained by the model. The taxonomic framework is in many cases classified based on Meunier and Danelian (2022 and 2023) at the species level, and higher taxonomic ranks are classified, mainly based on Suzuki et al. (2021).

4.2.4 CNN training

Earlier studies that attempted to compare the accuracy of several CNNs on microfossil assemblages found that Resnet50 usually performed very well for this purpose (Marchant et al., 2020; Tetard et al., 2020; Mitra et al., 2018); we therefore chose to apply this model for this particular study, instead of MobileNet v1 (Howard et al., 2017), which was we used previously (Carlsson et al. 2022) in study focusing on eight closely related species of the Eocene genus Podocyrtis. ResNet50 is a deep Convolutional Neural Network architecture (He et al., 2016), and is one of the variations of the ResNet (short for "Residual Network") family of models. The idea behind the formation of ResNet50 is to use residual learning to avoid disappearing gradients in very deep neural networks. Because when the networks get deeper, it becomes more difficult to update the weights of the earlier layers through backpropagation and by using residual learning; the network can thus propagate the gradient signal more easily, which improves the training of deeper networks. The weights of ResNet50 have also been pre-trained on a large dataset, for instance ImageNet, which includes millions of labelled images of about 1,000 classes (He et al., 2016). Our training set consisted of 80 % randomized images, chosen for each individual class present in the database, while the remaining 20 % was used for validation.

4.2.5 Performance metrics

The CNN training calculates automatically the classification accuracy and recall values based on the labeled dataset; both of them represent different ways of displaying the CNN performance (Fig. 4.2). For instance, if the aim is to detect all specimens belonging to a specific species (high recall value), sacrificing accuracy by including other objects not belonging to that species might be acceptable. This allows for easier tracking of the true abundance of that
specific species, and misclassified objects may be identified and ignored. However, if misclassified objects, which may belong to another important species, end up in another species class, this would lower their recall value. Therefore, it is important to have a high accuracy overall, but when examining individual classes separately, recall value holds a significant importance. Both indices are of great interest for different applications, more focused on biostratigraphy or paleoceanography for example.



Figure 4.2. A theoretical example showing the importance of a high accuracy and a high recall value for individual classes.

4.2.6. Test set to validate the CNN

To validate the consistency of the CNN training and testing from our dataset, we once again estimated how accurately the trained CNN performed and we compared it with a human operator. The neural network training produced a prediction model that was inserted into ParticleTrieur version 3.0.0. A threshold value can be set directly in ParticleTrieur before classifying new images. The threshold value constrains the degree of accuracy desired for an image needed to be classified into a given class. If the probability for an image/specimen to be classified to a specific class is too low, this image will be left unclassified. We randomly let the ParticleTrieur software pick 200 images selected from four new samples (800 in total), which were unlabeled and contained all particles segmented from composed images, from ODP site 1260, coming from different intervals within those used initially to train the CNN, which were 1260A-6R-3W, 18-20 cm; 1260A-8R-5W, 70-72 cm; 1260A-13R-5W, 66-68 cm; and 1260A-15R-4W, 69-71 cm. In ParticleTrieur we let the trained CNN identify all of them, using a threshold value of 0.5. We selected this low threshold value since it is better having more images classified, even if that will give a somewhat lower accuracy, recall, precision to classify more images rather than that they unclassified. We then examined afterwards what was correctly or incorrectly classified.

4.2.7 Application on 39 species

With the same four samples, we then this time entered all segmented particles from the entire mosaic composed images covering most parts of the coverslips, resulting in a total of 15 932 images, which were automatically classified with the CNN. Here we focused on the targeted 39 classes representing the selected well-known species. We applied different threshold (1-0.5) values for the pre-trained network and checked how the CNN could recognize these 39 most common nassellarian species. The aim of this test was to try future potential applications.

4.3 Results

4.3.1 Training of the initial dataset

Our trained CNN obtained 86.6 % in overall accuracy, with 75.6 % in precision/accuracy (Fig. 4.2); the latter measures the ability to avoid false positives and corresponds to the number of specimens classified as a class and also belonging to that class, divided by all specimens classified to that class. Our CNN obtained 78 % of recall, which calculates the ability to detect

all correct classifications, as it corresponds to the number of specimens in a class that were correctly classified divided by the total number of specimens in that class. The training and validation iterations are given in Supplements 4.2 figure S1, which confirms that the data is neither overfit, nor underfit.



PREDICTED LABELS

Figure 4.3. Simplified confusion matrix, showing the classification between different classes, with a focus on nassellarian super families. The x-axis on the left shows the true classes while the right axis shows recall values; the y-axis at the bottom shows the predicted classes while the top shows the precision value.

Some of the classes that are visible in the confusion matrix (Fig. 4.2 or Fig. S2 in Supplements 4.2 for a more detailed confusion matrix) show a low individual score (diagonal numbers), mainly due to the low number of available specimens (see number next to the label name), or to a high degree of resemblance between closely related or similar looking species. Most of the 39 important species of nassellarians we focused on were classified with a high accuracy (Fig. S2 in Supplements 4.2).

Species	Training (#)	Test (#)	Precision	Recall	F1 score
Dendrospyris stylophora	35	7	0.71	0.71	0.71
Elaphospyris didiceros group	295	59	0.85	0.90	0.88
Liriospyris clathrata group	77	15	0.78	0.93	0.85
Dictyomitra parva	229	46	0.98	0.98	0.98
Dictyoprora amphora group	124	25	0.75	0.72	0.73
Dictyoprora mongolfieri	271	54	0.89	0.94	0.92
Rhopalosyringium auriculaleporis	104	21	0.84	0.76	0.80
Rhopalosyringium biaurata	24	5	0.83	1.00	0.91
Carpocanopsis ornata group	21	4	0.33	0.50	0.40
Stichopterygium microporum	64	13	0.93	1.00	0.96
Sethochytris triconiscus	22	4	1.00	1.00	1.00

Table 4.1. Precision, recall and F1 score for the 39 most important nassellarian species focused in this study.

Species	Training (#)	Test (#)	Precision	Recall	F1 score
Lithochytris vespertilio	20	4	1.00	1.00	1.00
Lychnocanoma bajunensis	103	21	1.00	1.00	1.00
Lychnocanium babylonis group	56	11	0.75	0.82	0.78
Lophophaena radians group	48	10	0.70	0.70	0.70
Rhabdolithis pipa	28	6	0.75	1.00	0.86
Aphetocyrtis zamenhofi	40	8	0.60	0.75	0.67
Lophocyrtis alauda	17	3	0.75	1.00	0.86
Theocyrtis scolopax	22	4	0.00	0.00	0.00
Phormocyrtis embolum	130	26	0.71	0.92	0.80
Phormocyrtis lazari	68	14	0.77	0.71	0.74
Calocyclas hispida	29	6	0.57	0.67	
Podocyrtis (Lampterium) chalara	207	41	0.98	1.00	0.99
Podocyrtis (Lampterium) goetheana	115	23	1.00	1.00	1.00
Podocyrtis (Lampterium) mitra	184	37	0.97	1.00	0.99
Podocyrtis (Lampterium) sinuosa	93	19	1.00	0.79	0.88
Podocyrtis (Podocyrtoges) ampla	42	8	1.00	0.88	0.93
Podocyrtis (Podocyrtoges) diamesa	62	12	0.58	0.92	0.71
Podocyrtis (Podocyrtoges) phyxis	44	9	0.89	0.89	0.89
Podocyrtis (Podocyrtis) papalis	302	60	0.95	0.83	0.83

Species	Training (#)	Test (#)	Precision	Recall	F1 score
Podocyrtis (Lampterium) puellasinensis	24	5	0.80	0.80	0.80
Pentalocorys triacantha	100	20	0.78	0.90	0.84
Thyrsocyrtis rhizodon	43	9	1.00	0.67	0.80
Eusyringium fistuligerum group	38	8	0.86	0.75	0.80
Eusyringium lagena	16	3	0.50	0.67	0.57
Rhopalocanium ornatum group	21	4	0.67	0.67	0.57
Zealithapium mitra	33	7	0.54	1.00	0.70
Anthocyrtis mespilus group	53	11	0.78	0.64	0.70
Cycladophora spatiosa group	53	11	1.00	0.91	0.95

4.3.2 Performance validation from the test set

By using a threshold value of 0.5 the CNN classification resulted in that 769 images, out of a total of 800, were correctly classified. All classes were individually examined and the precision and recall values were calculated for each detected class (see Supplements 4.4). The CNN could classify these images into 76 different classes, while the human classifier considered that these groups belonged to 63 classes, excluding rare species which could not be classified into a proper class and were therefore referred to the class "others". Finally, an overall accuracy, recall, precision and F1 score were calculated (see Fig. 4.4).



Figure 4.4. CNN performance metrics based on a test set consisting of a total of 800 images, which were validated by a human expert.

4.3.3 Application on new samples

The manual validation of the classification made by a CNN on the four new samples using a threshold value of 1, always provided a correct interpretation without having any misclassified species, although it was far from detecting all classes and all true specimens for each class. Interestingly, with a threshold value of 0.9, we could in some sense detect almost all classes present in the dataset with a 69-73 % accuracy (see Supplements 4.5) and get some estimates for the number of taxa present in the samples. Some possible misclassifications can be easily reviewed at a later stage. Lowering the threshold values increased indeed the number of truly correct specimens in the right species, but it also increased the number of false classifications (see Supplements 4.5). Regarding the average score of predicted key-species using several threshold values (Table 4.2), the CNN was usually able to correctly identify between 24 to 34 out of the 39 targeted species and also not falsely detect species which do not exist (see Supplements 4.5).

 Table 4.2. Average prediction results of the 39 key-species and its accuracy for different threshold values along with the total number of images.

Threshold value	Predicted key- species	Correctly predicted species	key	Accuracy key- species	Predicted images	Total amount of images
1.00	77	77		1.00	1664	15,932
0.90	1144	784		0.69	11,074	15,932
0.80	1336	865		0.65	12,438	15,932
0.70	1501	943		0.63	13,462	15,932
0.60	1653	998		0.60	14,370	15,932
0.50	1799	1043		0.58	15,232	15,932

4.4 Discussion

4.4.1 Classification

Due to the high radiolarian diversity preserved in the studied samples obtained from the equatorial Atlantic, the correct assignment of every single image to a class may be challenging. One of the particularities of the studied radiolarian fauna is that it contains a lot of rare and undescribed species. In addition, the current state of the art does not allow confident taxonomic divisions in higher classes, as there are often unclear taxonomic boundaries. A lot of taxa are also morphologically similar to each other, and a lot of similarities sometimes exist between different families, giving us often hard time to find for them a proper affinity and even acceptable taxonomic names. As an example, we may state the confusion of *Carpocanopsis ornata* group with juvenile/broken Nassellaria group B and Pterocorythoid group (see the catalogue in Supplements 4.3 and the confusion matrix in Fig. S2, Supplements 4.2).

There have been challenges in determining the most commonly occurring nassellarian species. We aimed to ensure taxonomic clarity within these groups and only included taxa for which we are very confident for their assignment to that specific class. Consequently, numerous similar-looking species have been excluded from the classified dataset. This is because they did not fit into other groups or we judged that they would be confused with the taxon they most resemble with. However, this may be challenging as we strive to represent as much as possible from these samples, while also collecting new samples for automatic classification. As is the case of an example mentioned above, we encountered many specimens that resemble *Carpocanopsis ornata* group, but which do not always display the discriminatory characters of the species. Those specimens would also always be misidentified as *C. ornata* group and were therefore removed from the dataset. It is possible that these forms may be misidentified as *C. ornata* group in future samples, but they have to be looked upon individually afterwards. The important point is to obtain a dataset that consists of clean *C. ornata* group specimens identified with high confidence, so that there can be a clear reference of what a *C. ornata* group looks like (see the catalogue in Supplements 4.3 and the confusion matrix in Fig. S2, Supplements 4.2).

In addition, difficulties have also been met when we attempted to consistently classify a high number of images. Also, differences in individual specimen orientation and bubble inclusions played a big role in getting the CNN to work and to find the proper classes.

Since all objects appearing on a slide are trained and given a class, some of the included classes may be artificially defined and therefore correspond to taxonomically "bin" classes. We focused mainly on nassellarian radiolarians, trying to include as many classes as possible neatly defined at the species level. Although, for some nassellarian classes presented at higher taxonomic levels, we were obliged to accept a very large taxonomic concept. Regarding spumellarians, as the recognition of their inner structures is important, but difficult to detect with computer vision, identifications are even more challenging.

4.4.2 CNN training and new test set score

It is not a shock that the test accuracy, which randomly selected 800 images from four new samples, performed less than the test of the 20 % of the labeled dataset, from which 80 % were used to train the CNN. The accuracy of the test is 75.69 %, whereas the training accuracy is 85.6 %. From our labeled dataset, we have purposely removed a lot of "trash" particles, that will say particles which are broken fragments of radiolarians, blurry background particles etc. because if we kept these images, the CNN would rather be overtrained by the thousands of trash images and perform less. Besides that, rare species which consisted of too few species were left untrained by the neural network and therefore the 20 % of test did not include that many "trash" images or any rare species, in contrast to the new test set, which were completely randomly selected among any kind of particle that had been segmented. This last test was just to confirm how well the CNN was generally trained. For our last application we tried to just focus on the 39 well known radiolarian species, since it is the radiolarians that are of interest.

4.4.3 Feedback on individual species

By examining the 39 targeted species individually in every single sample, we observed that some species were easily identified correctly, while others performed poorly during the CNN training iteration and ended up in different classes. Table 4.3 presents the examined species and samples, along with their training performance. This arrangement facilitates a better understanding of the high or low number of correctly predicted species based on the training performance of the CNN.

Table 4.3. Training accuracy, recall and number of correctly identified specimens for each one of the 39 targeted species in this study and sample using a threshold value of 0.5, which is the lowest value we used in the test to identify all species which have an identification correctness equal or higher than 0.5. Hyphens "-" correspond to species not found in the samples (Meunier and Danelian, 2022; Meunier and Danelian, 2023).

Species	Accuracy	Recall	1260A 6R-3 W, 18-20 cm	1260A 8R-5 W, 70-72 cm	1260A 13R- 5 W, 66- 68 cm	1260A 15R- 4 W, 69- 71 cm
Dendrospyris stylophora	0.71	0.71	1	4	1	0
Elaphospyris didiceros group	0.85	0.90	16	27	10	48
Liriospyris clathrata group	0.78	0.93	12	2	7	10
Dictyomitra parva	0.98	0.98	_	_	_	108
Dictyoprora mongolfieri	0.89	0.94	36	9	38	14
Dictyoprora amphora group	0.75	0.72	1	2	8	68
Rhopalosyringium? auriculaleporis	0.84	0.76	_	2	3	3
Rhopalosyringium? biaurata	0.83	1.00	_	0	1	0
Carpocanopsis ornata group	0.33	0.50	12	0	1	_
Stichopterygium microporum	0.93	1.00	2	3	6	3
Sethochytris triconiscus	1.00	1.00	0	10	-	_
Lithochytris vespertilio	1.00	1.00	_	3	3	6
Lychnocanoma bajunensis	1.00	1.00	46	25	23	10
Lychnocanium babylonis group	0.75	0.82	_	5	14	8
Lophophaena radians group	0.70	0.70	2	0	2	_

Species	Accuracy	Recall	1260A 6R-3 W, 18-20 cm	1260A 8R-5 W, 70-72 cm	1260A 13R- 5 W, 66- 68 cm	1260A 15R- 4 W, 69- 71 cm	
Rhabdolithis pipa	0.75	1.00	2	6	1	3	
Lophocyrtis alauda	0.75	1.00	-	-	8	9	
Aphetocyrtis zamenhofi	0.60	0.75	_	_	0	20	
Theocyrtis scolopax	0.00	0.00	_	_	2	1	
Calocyclas hispida	0.57	0.67	4	12	0	0	
Phormocyrtis embolum	0.71	0.92	_	_	16	3	
Phormocyrtis lazari	0.77	0.71	9	4	_	_	
Podocyrtis (Lampterium) chalara	0.98	1.00 7 20		20	_	_	
Podocyrtis (Lampterium) goetheana	1.00	1.00	2	_	_	_	
Podocyrtis (Lampterium) mitra	0.97	1.00	_	2	1	_	
Podocyrtis (Lampterium) sinuosa	1.00	0.79	_	_	1	2	
Podocyrtis (Podocyrtoges) ampla	1.00	0.88	_	_	4	0	
Podocyrtis (Podocyrtoges) phyxis	0.89	0.89	_	_	-	0	
Podocyrtis (Podocyrtoges) diamesa	0.58	0.92	_	_	-	0	
Podocyrtis (Podocyrtis) papalis	0.94	0.83	8	6	1	6	
Podocyrtis (Lampterium) puellasinensis	0.80	0.80	0	0	_	_	
Thyrsocyrtis rhizodon	1.00	0.67	15	2	4	11	
Pentalocorys triacantha	0.78	0.90	12	43	7	4	

Species	Accuracy	Recall	1260A 6R-3 W, 18-20 cm	1260A 8R-5 W, 70-72 cm	1260A 13R- 5 W, 66- 68 cm	1260A 15R- 4 W, 69- 71 cm
Eusyringium lagena	0.50	0.67	_	_	3	7
Eusyringium fistuligerum group	0.86	0.75	9	5	2	0
Rhopalocanium ornatum group	0.67	0.67	3	0	6	1
Cycladophora spatiosa group	1.00	0.91	31	20	2	1
Anthocyrtis mespilus group	0.78	0.64	14	3	10	3
Zealithapium mitra	0.54	1.00	7	7	7	1

In general, the classes with good performance are *Elaphospyris didiceros* group (Plate 4.1.B), *Dictyomitra parva* (Plate 4.1.J), *Sethochytris triconiscus* (Plate 4.1.L), *Lithochytris vespertilio* (Plate 4.1.K), *Lychnocanoma bajunensis* (Plate 4.1.N), *Lychnocanium babylonis* group (Plate 4.1.M) and *Thyrsocyrts* (*Pentalocorys*) *triacantha* (Plate 4.1.Q). They are well-classified with few misclassified objects in their respective classes and they rarely appear in other classes.

The CNN was able to detect some of the true specimens of *Dendrospyris stylophora* (Plate 4.1.A). However, in many cases some trissocyclids/cephalospyrids with long feet were also misclassified as *D. stylophora*. Since this class is quite rare, it is difficult to make any accurate estimate about the application accuracy. For the training iteration the CNN obtained an accuracy of about 70 %.

Liriospyris clathrata group is a simple single-segmented nassellarian with large pores on its cephalis (Plate 4.1.C). Occasionally, some specimens may be misclassified into higher-ranked taxonomic classes, but overall, it performs well. It has about 80 % accuracy in the CNN training.

Overall, the CNN was able to correctly identify all forms of *Dictyoprora amphora* (Plate 4.1.G) group, but many other broken and unusual radiolarian forms, including some *Dictyoprora* spp. and *Dictyoprora mongolfieri* (Plate 4.1.F), were confused with this species group. Although the majority *of D. mongolfieri*) were identified in their true class, a smaller number of radiolarians or objects in the class were misclassified, indicating that these particular classes have been well-trained with a 90 % accuracy.

Although some *Rhopalosyringium* ? *auriculaleporis* specimens (Plate 4.1.H) were identified correctly, there are still several other objects that are misclassified as this species. The same goes for *R*. ? *biauratum* (Plate 4.1.I), but since not many have been observed in our test samples, it is difficult to make any fair estimates for this particular species.

Carpocanopsis ornata group (Plate 4.1.P) is poorly trained and is largely misclassified in the training process, with a training accuracy of only ca. 30 %. This is likely due to its very simple, smooth outline that is similar to many other taxa.

Stichopterigyum microporum (Plate 4.1.O) is classified well, despite having many similarities with *Euctyrtidium levisaltarix*, a species that was not individualized in this study, but was included in the *Eucyrtidium* genus class. Occasionally, some of these species may be mixed up if there are no morphologically distinct morphotypes. However, in those cases where there are distinct morphotypes, they are classified correctly.

The training accuracy of *Lophophaena radians* group (Plate 4.1.E) was ca. 70%, although there are not many estimates on how well this species is classified in the new samples obtained.

Rhabdolithis pipa (Plate 4.1.D) was detected frequently in our samples. However, a lot of other particles also ended up being included in this class, alike some types of spumellarians, since *R*. *pipa* has only one simple segment and it does not display any radial symmetry and bears two very long spines.

Lophocyrtis alauda (Plate 4.2.A) is well detected in samples coming from 1260A-13R-5W, 66-68 cm and 1260A-15R-4W, 69-71 cm. The somewhat similar looking species, *Aphetocyrtis zamenhofi* (Meunier and Danelian 2023) was only found present in 1260A-15R-4W, 69-71 cm but other objects or specimens were also mistakenly classified as *A. zamenhofi*, even at samples in which they do not exist. The training accuracies are 60 % for *A. zamenhofi* and 75 % for *L. alauda*.

Not many specimens of *Calocyclas hispida* (Plate 4.2.O) have been trained by the CNN, and it is only in sample 1260A-8R-5W, 70-72 cm that they appeared more often; they were classified well, without having a lot of misidentified radiolaria or other objects appearing in that class.

We did not encounter many typical specimens of *Podocyrtis goetheana* (Plate 4.2.F), but mainly early/transitional forms that do not display the elongated abdomen with the typical long straight bars on the median row of pores. In any case, most transitional forms were classified as *Podocyrtis chalara* (Plate 4.2.G) and exceptionally as *P. goetheana*. Otherwise, *P. goetheana* has a unique morphology compared to the rest of radiolarians present in our samples and it was therefore trained very well with a perfect F1 score; both accuracy and precision were 100 %.

Podocyrtis chalara (Plate 4.2.G) is well classified and recognized by the CNN and is easy to detect in our test samples. However, when it comes to *Podocyrtis mitra* (Plate 4.2.H), our

samples contained transitional forms between *Podocyrtis sinuosa* (Plate 4.2.I) and *P. mitra* or *P. mitra* and *P. chalara*. In the latter case, most specimens we captured were actually closer to *P. chalara* than *P. mitra* and were therefore classified as *P. chalara* rather than *P. mitra*. In the studied material there were too few typical specimens of *P. sinuosa* and *P. mitra* to make up a clear mind, but most of them were transitional forms and the CNN had two specimens classified as *P. sinuosa* and one as *P. mitra*. Individually, *P. sinuosa* from other samples are rather well-detected, even though it happened to have samples with a lower abundance of *P. sinuosa*.

The CNN can detect well *Podocyrtis papalis* (Plate 4.2.J). However, some forms that do not belong to *P. papalis* were incorrectly classified, quite often as *Podocyrtis ampla* (Plate 4.2.K) or *Podocyrtis diamesa* (Plate 4.2.M). *Theocyrtis scolopax* was also found classified as *P. papalis*. There were not many specimens of *P. diamesa* in these samples; therefore, none was classified as *P. diamesa*, and the few specimens of *P. diamesa* were actually classified as *P. papalis*, which is logical since they are very similar (see also Carlsson et al., 2022). Finally, *P. ampla* was always confused with *P. papalis*, as regrettably the CNN could not correctly detect any single *P. ampla*.

Podocyrtis phyxis (Plate 4.2.L), an important biostratigraphic index species occurring only within a short interval, was trained in the CNN with an 89 % accuracy and a recall value of 88 %; however, it was never classified correctly into its own class in the new test set of four samples. Instead, it was frequently misidentified as *Thyrsocyrtis rhizodon* (Plate 4.2.P), which is understandable, given that both species have an equal number of segments and consist of a horn and feet (although they differ in size and shape) and are more or less barrel shaped.

Unfortunately, the CNN cannot detect size differences because all images are resized to the same dimensions.

In conclusion, the *Podocyrtoges* lineage, which includes *P. ampla*, *P. phyxis*, and *P. diamesa*, cannot be reliably detected in the new set of test samples using our currently trained CNN. The different morphospecies of this lineage are difficult to be identified correctly due to the frequent occurrence of transitional forms that look very similar to other taxa in our dataset. Although we have a sufficient dataset of these species, more data and adjustments to similar-looking classes may be necessary to allow the CNN to more clearly distinguish them with a high degree of accuracy, as humans are able to do.

We did not encounter any *Podocyrtis puellasinensis* (Plate 4.2.N) species but it was originally trained with an accuracy and precision of 80 %.

Eusyringium lagena (Plate 4.2.R) can be easily detected by the CNN, and *Eusyringium fistuligerum* group (Plate 4.2.S) is occasionally misclassified as *L. vespertilio* (Plate. 4.1.K) or *S. triconiscus* (Plate 4.1.L). This is understandable since their proximal parts (thorax, cephalis, and thick conical horn) look almost identical.

Not many specimens of *Rhopalocanium ornatum* group (Plate 4.2.T) were detected. This species was trained on a small number of specimens and therefore only obtained a training accuracy of about 70 %. Some specimens of the *R. ornatum* group were found in its true class but were also appearing in other species and higher taxonomic leveled classes, which implies a lower recall number.

Both *Cycladophora spatiosa* group (Plate 4.1.Q) and *Anthocyrtis mespilus* group (Plate 4.1.R) could be detected well with a high accuracy in the samples in which they existed. In other samples, they had a lower prediction accuracy with more specimens incorrectly classified as either *A. spatiosa* group or *A. mespilus* group.

The classification of *Zealithapium mitra* (Plate 4.1.S) is not reliable due to its insufficient training dataset, which comprises only a small number of images. As a result, many broken radiolarians with large pores are frequently misidentified as *Z. mitra*, despite the fact that this species is characterized by large pores with a more conical shape.

The results we obtained suggest that many of the classes we distinguished may be confidently used in future applications (biostratigraphic or paleoceanographic). Well distinct species that the CNN can easily detect in whole assemblage analyses have a low recall value and appear rarely in wrong classes. As in many cases, the presence/absence of an index species is sufficient for biostratigraphy, the automated classification of targeted species in whole assemblage studies described in this paper, enables us to quickly confirm the correct classification of species and thus opens new perspectives for the application of Artificial Intelligence to radiolarian biostratigraphic studies. Apart from the 39 targeted species, many of the other classes had a worse performance; indeed, many half-complete or blurry specimens were classified as other objects. This makes it difficult for the moment to fully trust the CNN classification for counting all radiolarian species in order to get information about their relative abundances.

4.5 Conclusions

The newly established dataset of middle Eocene tropical radiolarians is well adjusted to fit a CNN. We obtained a high training accuracy of 86.6 % for its training in a CNN.

We evaluated the performance of our trained Convolutional Neural Network (CNN) on new tests and compared it to human performance, and obtained a testing accuracy of about 75.69 %. We thereafter specifically focused on 39 different species which the CNN demonstrated notable success in accurately identifying those species that had been well-trained.

In order to obtain an acceptable accuracy of the CNN for further studies, the labeling of classes was also revised to groups or separate taxa and reached the best compromise between CNN accuracy and consistent taxonomy. For example, when two morphologically close species or subspecies where often confused by the CNN, we found it better to fuse them together in an acceptable taxonomic framework, unless they were individually of biostratigraphic importance, instead of artificially biasing the CNN accuracy by often mistaking these two taxa with each other, in the same way as they may be confused by an operator.

This has proved to be an efficient way, both in speed and easiness, to quickly see what kind of radiolarian species exist and how many of them. However, since we focused only on a few classes, we cannot compare the relative abundance with certain taxa in relation to all radiolarians yet, but with improved methods and building a stronger dataset, it will be possible to get a closer estimate of the relative abundance of many taxa. This also highlights the importance of building good taxonomic datasets.

Overall, applying automatic image classification to the studied samples is time-saving, particularly for detecting the presence of the selected nassellarian species. This approach eliminates the need to manually count and track by an operator the targeted taxa present in a sample and avoids the risk of identification bias between different operators.

Data availability

Microscope slides from Leg 207, Hole 1258A, 1259A and 1260A, which were used for our dataset and application to a trained CNN, are stored at the University of Lille, France. The dataset (<u>https://doi.org/10.57745/E9YXW6</u>, Carlsson, 2023) is published in the University of Lille repository at Recherche Data Gouv.

Acknowledgments

This study was partly funded by the French government through the program "Investissements d'avenir" (I-ULNE SITE/ANR-16-IDEX-0004 ULNE) managed by the National Research Agency. It also received funding from the European Union's Horizon 2020 research and innovation program under the Marie Sklodowska-Curie grant agreement no. 847568. It was also supported by UMR 8198 Evo-Eco-Paléo and IRCICA (CNRS and Univ. Lille USR-3380).

A special thanks to Sylvie Regnier for helping to prepare the new samples by removal of nonsiliceous materials, and also many thanks to Ross Marchant for the quick reparation of bugs in ParticleTrieur.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <u>https://doi.org/10.1016/j.marmicro.2023.102268</u>.

<u>CHAPTER 4</u>: Subchapter 2 - Initial middle Eocene radiolarians dataset from the tropical Atlantic (ODP Leg 207) partly classified by K-means clustering

4.1 (2) Introduction

In our initial dataset, we attempted to classify all other particles or less interesting or resisting radiolarians by using unsupervised methods, like K-means clustering (Lloyd, 1957; MacQueen, 1967). K-means clustering is an unsupervised learning algorithm that classifies objects into a specified number of clusters (k) depending on similar data points. The idea is to divide a set of data points into k clusters, in such a way that the sum of the squared distances between each data point and its nearest cluster center is minimized. This is accomplished by an iterative process. In this process, each data point is assigned to the closest cluster center, and the center is then updated to reflect the average position of the data points assigned to it. This continues iteratively until no more changes occur in the centroids or until the iterations are complete. For images, each pixel corresponds to a data point and the k-means clustering algorithm tries to cluster similar pixels together.

The main difference between KNN and K-means clustering is that KNN is supervised while Kmeans clustering is unsupervised. KNN predicts the class label of a new data point which is based on the majority label of its k-nearest neighbors, while K-means clustering partitions a set of data points into k clusters based on their similarity. As in Subchapter 1, this dataset was trained in a ResNet50 (He et al., 2016) Convolutional Neural Network (CNN).

4.2 (2) Results

The result of this was not very satisfactory, especially since the majority of the images were not of interest and could therefore be overtrained and destroy the result of the classes that we are interested in. The overall accuracy for the first trial was about 61.4%, with a precision of 55.1% and 63.9% in recall.

4.3 (2) Discussion

In the first trial, we included a total of 163 classes, some of which were classified in a supervised way, such as 79 classes of Nassellaria and 7 classes of Spumellaria, all of which were double-checked by two other radiolarian taxonomists.

Our focus was mainly on classifying Nassellaria to an approved taxonomic rank, with a focus on 37 taxa at the species level, some of which are important in biostratigraphy. We tried to have at least 30 images per class, although some very abundant classes had a few hundred images in their class. Even though we were sure about a class sometimes, due to their low abundance, they were placed into higher taxonomic ranks. For the rest of Nassellaria, they were placed in higher ranks, for which we were comfortable with their identification, either at genus, family, superfamily, or even order level; some classes were even considered as "bin classes".

The remaining Spumellaria, together with other objects, were classified using unsupervised Kmeans clustering. The reason why Spumellaria are not more finely classified is twofold. The first reason is that Spumellarian identification is dependent on the morphology of the inner shell, which is not very visible when working with stacked images; the second reason is time constraints. However, since most Spumellaria have radial symmetry, they work better for K- means clustering, as clusters are also dependent on orientation and many Spumellaria look identical at any orientation they appear in.

By training the entire dataset with the high number of unsupervised classified "bin" classes, the CNN did not perform that well, since the majority of the dataset consisted of "bin" classes rather than the radiolaria of interest, which could create an imbalance in the dataset. The network might focus more on the "bin" classes due to their higher representation, ignoring the essential classes, which could affect the overall performance of the model. However, after revision based on the first results, changes in the dataset were made by removing a lot of the trash classes which are not of interest removing about 25,000 images, and just kept a few classes from those "bin" classes consisting of other fragments. The Spumellarians were better organized, although their taxa name may not be perfectly correct. For the Nassellaria, some higher taxonomic ranks consisted of too many variations in their datasets, which caused some problems. Therefore, many Nassellaria classes were divided further into lower-ranked classes, we still kept most of the species classes just gained a few more. This time we tried to obtain at least ten specimens per class instead of the previous 30 specimens per class. Taxa classes including very few numbers of specimens that were often misclassified as another class were often fused forming a new class and vice versa, classes with too much variation were split into several more classes.

Apart from this study, we initially tried to classify everything, both using K-means clustering and Spiking Time Dependent Plasticity (STDP) (Masquelier and Thorpe, 2007). However, we did not properly succeed due to a lot of irregularities in the images, such as differences in contrasts, touching objects, bubble inclusion (giving a lot of unnecessary pixel values), and Nassellaria being randomly rotated and having a bilateral symmetry. This attempt did not improve the sorting of images, but it could work to group Spumellaria to some degree. Spiking Time Dependent Plasticity (STDP) is a new unsupervised learning algorithm that is still in its early research stage. It is inspired by the biological brain and focuses on the strength of the different neural signals.

Appendix A: Initial list of all the classes for the radiolarian image database from the Demerara Rise at the Tropical Atlantic Ocean.

ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Gen indet sp A ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Gen indet sp B ACANTHODESMIOIDEA group B ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Dendrospyris stylophora ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Desmospyris obtusus group ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Elaphospyris didiceros grou ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE group A ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE group B ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE group C ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE group D ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Liriospyris clathrata group ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Petalospyris spp ACANTHODESMIOIDEA group A ACANTHODESMIOIDEA group C ACANTHODESMIOIDEA STEPHANIIDAE Zygocircus buetschli ACANTHODESMIOIDEA STEPHANIIDAE Zygocircus spp ARCHAEODICTYOMITROIDEA ARCHAEODICTYOMITRIDAE Dictyomitra parva ARTOSTROBIOIDEA ARTOSTROBIIDAE Dictyoprora amphora group ARTOSTROBIOIDEA ARTOSTROBIIDAE Dictyoprora mongolfieri ARTOSTROBIOIDEA ARTOSTROBIIDAE Dictyoprora spp ARTOSTROBIOIDEA ARTOSTROBIIDAE Siphocampe spp ARTOSTROBIOIDEA group A CARPOCANIOIDEA CARPOCANIIDAE Carpocanopsis ornata group CLADE G group A EUCYRTIDIOIDEA EUCYRTIDIIDAE Eucyrtidium spp EUCYRTIDIOIDEA EUCYRTIDIIDAE Stichopterygium microporum LITHOCHYTRIDOIDEA LITHOCYRTIDIDAE Lithochyrtis spp LITHOCHYTRIDOIDEA LITHOCYRTIDIDAE group A LITHOCHYTRIDOIDEA LITHOCYRTIDIDAE Lychnocanissa bajunensis LITHOCHYTRIDOIDEA LITHOCYRTIDIDAE Lychnocanium babylonis group NASSELLARIA group A NASSELLARIA group B NASSELLARIA group C NASSELLARIA group D NASSELLARIA group E NASSELLARIA MIXED KNOWN NASSELLARIA UNKNOWN Other Other (10)Other (11)Other (12)Other (13)

Other (14)

Other (15)
Other (16)
Other (17)
Other (18)
Other (19)
Other (2)
Other (20)
Other (21)
Other (22)
Other (23)
Other (24)
Other (25)
Other (26)
Other (27)
Other (28)
Other (29)
Other (3)
Other (30)
Other (31)
Other (32)
Other (33)
Other (34)
Other (35)
Other (36)
Other (37)
Other (38)
Other (39)
Other (4)
Other (40)
Other (41)
Other (42)
Other (43)
Other (44)
Other (45)
Other (46)
Other (47)
Other (48)
Other (49)
Other (5)
Other (50)
Other (6)
Other (7)
Other (8)
Other (9)
PLAGIACANTHOIDEA - DICTYOCRYPHALIDAE - Dictyocryphalus capito gro
PLAGIACANTHOIDEA CERATOCYRTIDAE Ceratocyrtis spp

PLAGIACANTHOIDEA DICTYOCRYPHALIDAE Dictyocryphalus radians grou PLAGIACANTHOIDEA group A PLAGIACANTHOIDEA LOPHOPHANIDAE group A PLAGIACANTHOIDEA LOPHOPHANIDAE group B PLAGIACANTHOIDEA PLAGIACANTHIDAE Rhabdolithis pipa PLECTOPYRAMIDOIDEA PLECTOPYRAMIDIEA group A PTEROCORYTHOIDEA LOPHOCYRTIDAE group A PTEROCORYTHOIDEA LOPHOCYRTIDAE group B PTEROCORYTHOIDEA group A PTEROCORYTHOIDEA group B PTEROCORYTHOIDEA LOPHOCYRTIDAE Aphetocyrtis zamenhofi PTEROCORYTHOIDEA LOPHOCYRTIDAE Aphetocyrtis alauda PTEROCORYTHOIDEA LOPHOCYRTIDAE Paralampterium scolopax PTEROCORYTHOIDEA PTEROCORYTHIDAE APHEtocyrtis spp

PTEROCORYTHOIDEA PTEROCORYTHIDAE Calocyclas hispida

PTEROCORYTHOIDEA PTEROCORYTHIDAE Calocycloma ampulla PTEROCORYTHOIDEA PTEROCORYTHIDAE Phormocyrtis embolum PTEROCORYTHOIDEA PTEROCORYTHIDAE Phormocyrtis group PTEROCORYTHOIDEA PTEROCORYTHIDAE Phormocyrtis lazari

PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Lampterium chalara PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Lampterium goetheana PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Lampterium mitra PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Lampterium sinuosa PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis papalis

PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Podocyrtoges ampla PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Podocyrtoges diamesa PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Podocyrtoges phyxis PTEROCORYTHOIDEA THEOCOTYLIDAE Axocorys spp PTEROCORYTHOIDEA THEOCOTYLIDAE group A PTEROCORYTHOIDEA THEOCOTYLIDAE Theocorys spp PTEROCORYTHOIDEA THEOCOTYLIDAE Thyrsocyrtis rhizodon group PTEROCORYTHOIDEA THEOCOTYLIDAE Pentalocorys triacantha group PTEROCORYTHOIDEA THEOPERIDAE Eusyringium fistuligerum group PTEROCORYTHOIDEA THEOPERIDAE Eusyringium lagena PTEROCORYTHOIDEA THEOPERIDAE group A PYLOBOTRYDOIDEA PYLOBOTRYDIDAE group A Spumellaria SPUMELLARIA STYLOSPHAEROIDEA - STYLATRACTIDAE Stylophaera spp Spumellaria (10) Spumellaria (14) Spumellaria (15) Spumellaria (17)

Spumellaria (18) Spumellaria (2) Spumellaria (21) Spumellaria (24) Spumellaria (26) Spumellaria (27) Spumellaria (28) Spumellaria (29) Spumellaria (3) Spumellaria (30) Spumellaria (5) Spumellaria (6) Spumellaria (9) SPUMELLARIA LITHOCYCLIOIDEA LITHOCYCLIIDAE Lithocyclia ocellus SPUMELLARIA LITHOCYCLIOIDEA PHACODISCIDAE Periphaena decora SPUMELLARIA LITHOCYCLOIDEA LITHOCYCLIIDAE Heliosestrum spp SPUMELLARIA PHORTICIOIDEA HISTIASTRIDAE Histiastrum group STYLOSPHAEROIDEA STYLATRACTIDAE Zealithapium mitra STYLOSPHAEROIDEA STYLATRACTIDAE Zealithapium spp THEOPILIOIDEA ANTHOCYRTIDIDAE Anthocyrtis group THEOPILIOIDEA ANTHOCYRTIDIDAE Anthocyrtis mespilus group THEOPILIOIDEA ANTHOCYRTIDIDAE Anthocyrtis spatiosa TREMATODISCOIDEA A TREMATODISCOIDEA B TREMATODISCOIDEA C TREMATODISCOIDEA D TREMATODISCOIDEA E TREMATODISCOIDEA F TREMATODISCOIDEA G TREMATODISCOIDEA H TREMATODISCOIDEA I

<u>CHAPTER 5</u>: Using convolutional neural network for improving the biostratigraphy on middle Eocene radiolaria from ODP Leg 207, Site 1260

Abstract

Radiolaria among other microfossils are useful to date sediments by studying the species occurences shifting between different intervals. Biostratigraphic work demands a lot of time and effort and the need for a taxonomic expert; therefore, automatic image recognition can simplify the workflow. Here, we worked with images from samples photographed by an automatic microscope and used an already trained CNN, to automatically classify 90,237 images coming from 23 different samples. Despite using far fewer samples than the original biostratigraphic work on the same ODP Leg 207 Site 1260 and a limited amount of sediment to avoid touching objects, so that the segmentation software can crop each particle into individual images, our results show that in most cases the species ranges found in this study compare well to the results of an earlier, classically conducted, biostratigraphic study. It is thus established that stratigraphic ranges may be well recognized by the CNN, and in some cases, the CNN could expand the ranges. However, due to the low number of sediments together with some more poorly trained species, the ranges of rare species cannot always be detected by the CNN. Furthermore, the biostratigraphic species trained were also revised, showing their potential usefulness in further studies. Some classes are very well-trained and could be used for automatic image recognition without any further taxonomic expertise, while other species with high recall values could be applied to semi-automatic image recognition for biostratigraphy

with some revision. A few classes still need more training before being reliably used for biostratigraphic applications.

5.1 Introduction

Radiolaria are microscopic unicellular heterotrophic eukaryotic plankton organisms, living exclusively in marine environments. It is made of aesthetically pleasing amorphous silica test or shells and has existed since the Cambrian age. Radiolaria are commonly used in biostratigraphy in areas where there are better studied calcareous microfossils such as foraminifera and nannofossils are absent.

Sediments from the Ocean Drilling Program (ODP) Leg 207, drilled at the Demerara Rise of the coast off Suriname in the tropical Atlantic Ocean; are enriched in silica resulting in highly preserved radiolarians. The middle Eocene sequence is almost entirely complete, with estimations of about 500 species of radiolarians.

ODP Site 1260 has been dated using both magnetostratigraphic and cyclostratigraphic methods (Westerhold and Röhl, 2013). A recent biostratigraphic work was obtained on this site by Meunier and Danelian (2022), in which the limits of radiolarian biozonations of RP16, RP15, RP14, and RP13 were determined, along with suggestions for new biozonations of RP15a and b; RP14a and b; and RP13a and b.

In recent years, a lot of studies have applied Artificial Intelligence (AI) to micropaleontology, especially using Convolutional Neural Networks (CNNs), which is a type of AI specifically designed for visual recognition of videos and images. Carlsson et al., (2023), trained a large image dataset from samples coming from ODP Leg 207, from Site 1258, 1259, and 1260 with

96 classes representing a large part of the entire middle Eocene assemblages from the tropical Atlantic, obtaining thus a training accuracy of about 86 %; and a testing accuracy of ca 76%.

This study is a continuation of Carlsson et al. (2023), by using a real biostratigraphic approach to test how a trained CNN works in identifying biostratigraphic events compared to another recent biostratigraphic work made by Meunier and Danelian (2022). Here, we only focused on biostratigraphic important species that have been trained by a CNN, using a Resnet50 architecture. We commented about the ranges of species that the CNN could detect which was manually supervised afterward and we also recognized which classes are more reliable than others for fully automatic or semi-automatic image recognition that eventually could be used for other taxonomic experts or non-taxonomists.

5.2 Materials and Methods

5.2.1 Materials

Deep sea sediment samples are retrieved from the Ocean Drilling Program (ODP) Leg 207, Site 1260, at the Demerara Rise, off the coast off Suriname. A total of 23 different middle Eocene samples (Table 5.1), with one sample containing eight small coverslips, which is 1.2 x 1.2 cm (Tetard et al., 2020), was used in this study.

EXPEDITION	SITE	HOLE	CORE	CORE TYPE	SECTION	TOP DEPTH	BOTTOM DEPTH	MBSF TOP	MCD TOP	AGE
207	1260	А	6	R	3	18	20	41.38	41.38	
207	1260	А	6	R	4	20	22	42.9	42.9	
207	1260	А	6	R	4	119	121	43.89	43.89	40,070516
207	1260	А	6	R	5	63	65	44.83	44.83	40,112968
207	1260	А	7	R	1	69	71	48.19	48.29	40,298900
207	1260	А	7	R	3	18	20	50.68	50.78	
207	1260	А	8	R	5	70	72	63.9	64	

 Table 5.1. List of samples used for the image recognition.

207	1260	А	8	R	6	65	67	65.35	65.45
207	1260	А	9	R	1	64	66	67.54	66.54
207	1260	А	9	R	2	64	66	69.04	68.04
207	1260	А	9	R	4	68	70	72.08	71.08
207	1260	А	10	R	4	64	66	81.74	80.84
207	1260	А	11	R	3	61	63	89.91	88.53
207	1260	А	11	R	7	68	70	95.98	94.6
207	1260	А	13	R	3	70	72	109.3	107.32
207	1260	А	13	R	4	68	70	110.78	108.8
207	1260	А	13	R	5	66	68	112.26	110.28
207	1260	А	13	R	6	68	70	113.78	111.8
207	1260	А	14	R	5	65	67	121.95	120.27
207	1260	А	15	R	1	69	71	125.59	123.78
207	1260	А	15	R	2	68	70	127.08	125.27
207	1260	А	15	R	3	69	71	128.59	126.78
207	1260	А	15	R	4	69	71	130.09	128.28

5.2.2 Methods

Samples were prepared using the same technique described in Tetard et al. (2020), and photographed in a Nikon automatic microscope, see Carlsson et al. (2023) for more details concerning the microscopy. Furthermore, the image processing is described in both Tetard et al. (2020) and Carlsson et al. (2023).

From all the 23 samples, 90,237 segmented images were recovered, and between 1,700 and 9,000 individual images for each sample. These were then analyzed by a trained CNN (Carlsson et al., 2023). The automatic image classification took part entirely in the software ParticleTrieur (Marchant et al., 2020) v.3.0.4. Images were pre-prepared to be equally in a 256-pixel size, from previously having 0.36 μ m/pixel resolution. A confidence score/threshold value of 0.3 was preselected, enabling more images to be classified, despite leading to a lot of images being classified incorrectly (Carlsson et al., 2023). The average running time for the image

classification in ParticleTrieur was about two images per second using a computer standard intel Core i5, 10th generation, hardware, taking less than two working days to finalize. A total of 90 078 images were classified and important biostratigraphic species were double-checked so that they ended up in the correct class.

5.3 Results

5.3.1 CNN prediction and revision

In ParticleTrieur, data such as sample counts were exported, showing species occurrence predicted by the CNN along with the samples, which is chronologically sorted starting at the youngest age. We have also removed all classes which are not known to be of biostratigraphic importance (see Table 5.2). Apart from that, images can be exported as well and are sorted at their predicted classes keeping the name of which sample they come from.

Table 5.2. Biostratigraphic important species that has been automatically recognized by a trained CNN from a batch of ca. 90 000 images covering 23 samples. The pale-green colour transition is the number of specimens/images recognized by the CNN, green being the highest.

	D. parva	R? auriculaleporis	R? biaurata	C. ornatum group	L. vespertilio	L. triconiscus	L. babylonis group	<i>L.a radians</i> group	A. zamenhofi	L. alauda	T. scolopax	P. embolum	P. lazari	P. L. chalara	P. goetheana	P. mitra group	P. sinuosa	P ampla	D diamoca	D nhuvic	E. lagena	Z. mitra
Sample																						
1260A-6R-3W-18-20	3	10	0	25	1	0	2	8	9	1	0	20	10	15	2	2	2	0	0	0	0	48
1260A-6R-4W-20-22	0	8	2	10	18	4	3	4	15	1	0	15	2	44	1 7	10	3	2	3	1	2	71
1260A-6R-4W-119-121	10	4	0	31	0	2	5	15	10	0	1	13	5	11	0	9	1	0	2	1	0	34
1260A-6R-5W-63-65	2	11	0	22	6	8	3	4	26	2	1	21	4	35	1 3	4	8	0	1	0	0	104
1260A-7R-1W-69-71	1	0	0	1	21	20	1	0	7	0	0	7	1	86	5	5	0	0	1	4	0	73
1260A-7R-3W-18-20	0	9	1	0	9	15	0	9	18	0	1	20	6	41	1	0	0	0	0	1	0	88

1260A-8R-5W-70-72	1	4	3	0	3	12	5	1	13	0	0	4	4	38	1	5	2	0	0	3	1	43
1260A-8R-6W-65-67	3	2	5	1	2	4	5	1	5	0	1	6	5	18	8	9	1	0	0	5	0	24
1260A-9R-1W-64-66	0	0	3	1	3	3	6	0	3	0	0	19	2	30	7	29	6	0	2	0	0	35
1260A-9R-2W-64-66	5	2	5	3	2	5	6	0	9	0	1	4	1	8	2	6	0	0	0	0	0	18
1260A-9R-4W-68-70	9	32	3	15	1	6	8	9	31	1	2	24	35	6	1	3	1	0	0	0	0	14
1260A-10R-4W-64-66	9	23	7	20	10	11	23	27	57	2	3	56	2	7	1	12	4	0	1	0	1	52
1260A-11R-3W-61-63	0	11	0	4	18	1	7	4	19	7	19	19	6	0	0	4	1	2	1	0	0	23
1260A-11R-7W-68-70	2	40	5	21	1	2	9	8	3	4	4	20	10	0	0	0	2	0	0	0	1	12
1260A-13R-3W-70-72	1	8	5	2	3	0	26	5	16	9	8	30	5	13	1	17	5	2	0	1	3	28
1260A-13R-4W-68-70	0	3	1	1	4	3	12	2	7	12	1	8	0	5	0	5	5	9	0	0	1	23
1260A-13R-5W-66-68	2	23	3	13	3	0	17	8	7	8	3	24	7	1	0	1	1	5	0	1	5	15
1260A-13R-6W-68-70	0	5	0	7	8	0	12	5	10	10	2	30	2	0	0	4	7	4	1	3	1	11
1260A-14R-5W-65-67	9	26	0	12	3	2	4	15	4	5	8	15	5	0	0	1	2	1	1	0	1	11
1260A-15R-1W-69-71	6	13	1	1	13	1	7	2	36	3	9	26	4	2	2	6	3	2	0	0	10	29
1260A-15R-2W-68-70	16	15	2	18	20	3	14	8	98	17	21	98	20	4	0	1	3	0	0	1	23	52
1260A-15R-3W-69-71	40	20	2	7	9	3	8	15	102	18	9	74	10	9	1	1	4	1	0	0	8	129
1260A-15R-4W-69-71	43	18	0	34	6	1	14	9	29	8	0	30	18	4	2	0	3	0	0	0	7	21

As we are using a high confidence value and a network that has an accuracy of about 76%, we do revise each class and remove the images from each biostratigraphic class that is misclassified, and thereafter we update the values again on the same type of table (see Table 5.3).

Table 5.3. Biostratigraphic important species having all the misclassified specimens removed from each class after revision. The yellow-green colour transition is the number of specimens/images recognized by the CNN, green being the highest. Gray names include species that was not originally trained by the CNN.

	D. parva	R? auriculaleporis	R? biaurata	C. ornatum group	L. vespertilio	L. triconiscus	L. babylonis group	L.a radians group	A. zamenhofi	L. alauda	T. scolopax	P. embolum	P. lazari	P. chalara	Podocyrtis sp. cf. P. goethea	P. goetheana	P. mitra	P. trachodes	P. sinuosa	P. ampla	P. diamesa	P. phyxis	E. lagena	Z. mitra
Sample															na									
1260A-6R-3W-18-20				11									10	3	4	1								7
1260A-6R-4W-20-22				6	16	3	4						1	23	11			1						19
1260A-6R-4W-119-121				10	0	0	0			1			4	6	1									4
1260A-6R-5W-63-65				4	5	8	3			0			4	21	8									44
1260A-7R-1W-69-71				0	19	14	1			0			1	67	1									17
1260A-7R-3W-18-20				0	8	14	0			0			6	31	0									5
1260A-8R-5W-70-72		2		0	3	10	5			0			4	24	1		2							9
1260A-8R-6W-65-67		0	3	0	2	4	3			0			2	11			0							8
1260A-9R-1W-64-66		0	0	0	3	3	6			0			0	14			1							3
1260A-9R-2W-64-66		1	1	0	2	5	5			0			1	4			1							6
1260A-9R-4W-68-70		10	2	4	1	6	6			0			34	0			1	1						2
1260A-10R-4W-64-66		1	3	0	9	10	14	3		0		47		0			2	9						
1260A-11R-3W-61-63		0	0	4	17		4	1		6	16	7		0				2		2				
1260A-11R-7W-68-70		2	0	2	1		9	0		3	2	15		0				2		0				
1260A-13R-3W-70-72		5	1		3		25	1		7	4	14		1					3	2			1	
1260A-13R-4W-68-70		1	0		3		10	1		11	1	6							1	9			0	
1260A-13R-5W-66-68		2	2		3		13	3		8	1	17							1	5			3	
1260A-13R-6W-68-70		0			8		9			10	0	5							7	4			1	
1260A-14R-5W-65-67		6			3		3		1	5	6	6							2	1			1	
1260A-15R-1W-69-71	3	5			13		3		20	1	8	14							2	2			8	
1260A-15R-2W-68-70	15	4			20		2		63	15	17	39							1			1	21	
1260A-15R-3W-69-71	37	7			9		5		53	12	4	41							3				8	
1260A-15R-4W-69-71	31	2			6		8		20	8		18							2				7	

5.3.2 Correlation with earlier biostratigraphic work

Our samples are retrieved from different intervals to complement for a higher resolution data. The results from Table 5.3 along with the earlier biostratigraphic work made on Site 1260 (Meunier and Danelian, 2022), were merged (see Fig. 5.1), displaying the different species occurences obtained from both Meunier and Danelian (2022) and from the CNN prediction in this study. Note that the taxonomic names from both studies vary, and to avoid any confusion see the taxonomic appendix.
New radiolarian subzones	Artobonys bia wita	Агговопу в ингиla leports Суртосагриин огнанин	Eusyringum lagena Lithochynis vesperalio	Lophocyrtis (Lophocyrtis ?) alauda	Lophophaena radians	Phormocyrtis embolum Podocyrtis (Lamptarium) chalara	Podocyriis (Lampterium) goetheana	Podocyrtis (Langtarium) sp. cf.P. (L.) goetheana Podocyrtis (Langtarium) mina	Podocyrtis (Lanpteriun) sinuo sa Podocyrtis (Lanpteriun) trachodes	Podo cyrtis (Podocyrtoges) ampla Podo cyrtis (Podocyrtoges) diamesa	Podocyrtis (Podocyrtoges) phyxis Sethochynis ? babylonis	Sethochypis? niconiscus	Siphocanye ? anygdala Theocristic secolonov	Zealithapium mitra
RP16														
RP15b														
RP15a														
RP15a/ RP14b				I										
RP14b														
RP14a														
RP13b								•						
RP13a														
RP12														

New ranges for this study

Ranges by Meunier and Danelian (2022)

Ranges by Meunier and Danelian (2022) not detected by the $\ensuremath{\mathsf{CNN}}$

Figure 5.1. Biostratigraphic ranges from Meunier and Danelian (2022) correlated with the ones obtained by the CNN in this study. Only the occurrence ranges from each species and the radiolarian zones including the sub biozones obtained in Meunier and Danelian (2022).

5.3.3 Biomarkers for the radiolarian zonation

RP16

No new limits of the base of RP16 based on the First Occurrence (FO) of *Podocyrtis* (*Lampterium*) goetheana (Haeckel, 1887). Only one true *P. goetheana* was found 1260A-6R-3W, which is from the youngest sample included in this study. However, based on a recent study by Pinto et al., (2023) a new definition of an intermediate form between *Podocyrtis* (*Lampterium*) chalara Riedel and Sanfilippo, 1970 and *P. goetheana* was suggested, which also affects the base of RP16. We found the intermediate form *Podocyrtis* sp. cf. *P. goetheana* (Plate 5.1.A) as far down as 8R-5W, 70-72 cm here, although the FO of the true *P. goetheana* suggested by Meunier and Danelian (2022) was first found at sample 1260A-6R-6W, 55-57 cm.

RP15b

The base of the newly suggested biozone RP15b (Meunier and Danelian, 2022) is based on the Last Occurrence (LO) of *Rhopalosyringium? biauratium* (Ehrenberg, 1875). This species was rare in our samples, and we did not get a lot of *R? biauratium*. Therefore, the LO of our samples are located at sample 1260A-8R-6W, 65-67 cm; while Meunier and Danelian (2022) recognized the LO at sample 1260A-8R-3W, 54-56 cm.

RP15/RP15a

The next biozone is RP15 or RP15a; and the base is recognized by the Evolutionary Transition (ET) between Podocyrtis (Lampterium) mitra Ehrenberg, 1854 and P. chalara. This limit is complicated in a lot of ways, the first one is that Podocyrtis (Lampterium) trachodes Riedel and Sanfilippo, 1970 and P. mitra are defined as the same species group by the trained CNN, secondly, we are tracking the limit of where the majority of the specimens for each species are located, and many forms apart from P. mitra also includes P. trachodes, which we separated manually for comparisons here. Then there are a lot of intermediate forms, and many taxonomists have different opinions on where to draw the line between different species and intermediates. In this case, we recognized that the base of RP15a is a bit diffused, we have recognized early forms of P. chalara rather than late forms of P. mitra, which is suspected to have occurred in the earlier study. The pores were carefully counted to have a maximum circumference of 10 or 12, and more than 13 is recognized as a P. mitra. Our base appears lower at sample 1260A-9R-2W, 64-66 instead of 1260A-8R-6W, 54-56, if only relying on the revised data from the CNN. However, the number of images from the CNN is not that many, therefore it is not correct to determine the base on RP15a from this study. It can also be very difficult for the CNN to distinguish between having more or less than 13 pores.

RP14b

The base of RP14b (Meunier and Danelian, 2022) is based in the FO *Sethochytris triconiscus* Haeckel, 1887. The FO in our samples is compared to the work made by Meunier and Danelian (2022) appearing a bit later at sample 10R-4W, 64-66 cm, probably due to the sample resolution, meaning that the RP14b limit remains from the FO found by Meunier and Danelian (2022) at sample 11R-2W, 55-57 cm.

RP14/RP14a

The ET between *Podocyrtis (Lampterium) sinuosa* Ehrenberg, 1875 and *P. mitra* marks the base of RP14 or RP14a. Since we are dealing with transitional forms and *P. mitra* again, things remain complicated and rare *Podocyrtis fasciolata* (Nigrini, 1974), could be mistaken as intermediate juvenile forms between *P. sinuosa* and *P. chalara*. It is worth taking in consideration that the feet in *P. fasciolata* are underdeveloped. All specimens interpreted to be *P. fasciolata* are removed from this study.

RP13b

The newly suggested biozone RP13b's base is defined by the LO of *Podocyrtis (Podocyrtoges) phyxis* Sanfilippo and Riedel, 1973, and since *P. phyxis* is poorly trained by the CNN, only one true specimen was recognized by the CNN, meaning that the previous limit made by Meunier and Danelian (2022) remains the same.

RP13/RP13a

The base of RP13 or RP13a is recognized by the ET between *P. physis* and *P. chalara*. Since both of these species are trained poorly resulting in a very low recall number, not a lot of specimens are recognized by the CNN and therefore the previous limit made by Meunier and Danelian (2022) remains the same for now.

5.3.4 Remaining bio events

Most of the bio events stay the same, and Figure 5.1 can be looked at in more detail for that. However, there are some differences worth mentioning and explaining.

The LO of *Lophocyrtis alauda* (Ehrenberg, 1875) was identified by Meunier and Danelian (2022), to occur at 1260A-10R-5W, 55-57, and the LO in this study identified the LO at 1260A-

11R-3W, 61-63, where they are more or less common. However, a single large specimen resembling *L. alauda* was found at sample 1260A-6R-5W, 63-65 (see Plate 5.1.B).

P. trachodes, is not trained by the CNN as an individual class and is instead included in *P. mitra* group, but to make this study comparable with the one obtained by Meunier and Danelian (2022), we therefore manually divided the specimen predicted as *P. mitra* group, into either *P. trachodes* or *P. mitra*, and we based our decision mainly on the rough outline of the abdomen. Still using a CNN saves a lot of time. Yet again, this class decision can be affected by different opinions between taxonomists. We placed the LO of *P. trachodes* at 1260A-6R-4W, 20-22, based on a single specimen (see Plate 5.1.C). *Podocyrtis mitra* is known to exist rarely around this interval.

Plate 5.1. A) An example of *Podocyrtis* sp. cf *P. goetheana* (Pinto et al., 2023) a typical intermediate form between *P. chalara* and *P. goetheana*. B) The single *A. alauda* specimen appearing at sample 1260A-6R-4W-119-121. C)
A single specimen of *P. trachodes* appearing at sample 1260A-6R-4W, 20-22.



5.4 Discussion

5.4.1 General discussion

Even from the fully trained CNN, one can sometimes see clusters of a larger number of identified specimens for certain species; already from there one may make approximate estimations on their appearances or where they might appear in a higher abundance. Clear examples of this can be seen within *Dictyomitra parva* (Kim,1992), *Rhopalosyringium? auriculaleporis* (Clark and Campbell, 1942), *L. triconsiscus*, *Aphetocyrtis zamenhofi* Meunier and Danelian, 2023, *Phormocyrtis embolum* Ehrenberg, 1875, *P. chalara* and *P. goetheana*.

Some species can be more reliable than others to use for an entirely automatic classification (high recall and high accuracy), while others can be used in a semi-automatic way (high recall), with or without trained taxonomists, which are depending on whether the morphospecies are easy to recognize by anyone. Some classes might also need to be trained or regrouped before being able to be used for any kind of automatic recognition. Below we revised each species class and their reliability.

The use of automatic image recognition can also remove human biases, such as the expectation of certain taxa, not appearing at different intervals, even though they may exist rarely but are not identified as that species, since the brain ignored or categorized the specimen as unknown.

5.4.2 Individual species revision and their credibility in future biostratigraphic work

A. zamenhofi

A lot of similar looking taxa are misclassified into this class, due to the low confidence score. However, most specimens belonging to *A. zamenhofi* is recognized here, having a high recall value, and this form is relatively easy to recognize. This class can be used for automatic image recognition by revision of non-taxonomists.

Carpocanopsis ornata (Ehrenberg, 1875) group

This class was trained with a 50 % accuracy in Carlsson et al., 2023; and was not expected to be performing well due to the similarities with a lot of other Nassellarian species. However, the CNN was still able to capture a lot of *C. ornata*. This class could both be trained better and be used in a semi-automatic way, with the revision of an Eocene taxonomist.

D. parva

Even though this material has not been used to calculate relative abundances, we can still see that there are a large number of *D. parva* in the older samples. There are still a lot of other particles mistaken in this class, but this species is very easy to recognize even by a non-taxonomist and therefore could be easily revised after automatic recognition by the CNN. In this study we did not find any *D. parva* in 14R-5W, 65-67 cm, a range they should exist in according to the biostratigraphic study by Meunier and Danelian (2022). One simple reason could be that they are not so abundant there.

Eusyringium lagena (Ehrenberg, 1875)

Very easy species to recognize and it has both a high accuracy and recall. This species class trained by this CNN could easily be applied for automatic image recognition without any major revisions.

L. vespertilio

Easily recognizable having both a high accuracy and a high recall value and therefore excellent to use for completely automatic recognition without the need for too much revision.

L. alauda

Relatively high recall value, and easy to be recognized for a trained taxonomist. For nontaxonomists, it could be difficult to separate specimens belonging to the Lophocyrtiidae family, since this family is very diverse. Therefore, this class trained by this CNN can be used for semiautomatic classification requiring a radiolarian specialist to revise.

Lophophaena radians Ehrenberg, 1875 group

L. radians were very rare in our samples, although many images/specimens were misclassified as *L. radians*. In the trained CNN, they have an accuracy and recall value of about 70 %, but since we use a low confidence score more images are expected to be misclassified.

Lychnocanium babylonis (Clark and Campbell, 1942) group

This group of species is well trained and therefore has a high accuracy and recall and could potentially be used for automatic image recognition without any major revisions.

P. embolum

P. embolum seems to be relatively common or abundant at the interval it exists in. It has a high recall value and is easy to recognize by a trained taxonomist.

Phormocyrtis lazari Meunier and Danelian, 2023

P. lazari has a high recall but only exists in the younger materials. The older materials are often mistaken for *Theocyrtis scolopax* (Ehrenberg, 1875), even though they do not exist during the same interval.

P. chalara

We have done no relative abundance estimations here therefore we cannot know what percentages of abundance the *P. chalara* have in relation to all radiolaria in specific samples.

However, *P. chalara* is one of the most abundant species among all these species, especially at the interval it exists. It has a high recall and accuracy and could be used for automatic classification. However, early *P. chalara* are even for taxonomist difficult to classify, but typical forms can be identified even by non-specialists.

Podocyrtis (Lampterium) goetheana (Haeckel, 1887)

P. goetheana is well-trained and has both a high accuracy and recall value. However, the samples from site 1260A, are not young enough to display the typical forms of *P. goetheana* with the very long second abdominal pore and long straight bars. This is one of the easiest radiolarian species to recognize even by a non-taxonomist.

P. mitra group

Not a lot of typical *P. mitra* were captured here, it can have to do with low abundance or low recall value, or their near transition from *P. sinuosa* to *P. mitra* and from *P. mitra* to *P. chalara*. Although, *P. trachodes* included in this group is always recognized.

P. sinuosa

This species class has a low recall value. It is known to be misclassified for *Podocyrtis papalis* Ehrenberg, 1847. Therefore, we do not obtain a lot of *P. sinuosa*.

Podocyrtis (Podocyrtoges) ampla Ehrenberg, 1875

P. ampla also has a low recall value and is often mistaken as other closely related *Podocyrtis* species which also have bad recall and accuracy values. However, the accuracy is high for *P. ampla*, as not a lot of other species are misclassified into this group.

Podocyrtis (Podocyrtoges) diamesa Riedel and Sanfilippo, 1970

No *P. diamesa* was recognized correctly by the CNN. Therefore, this trained CNN, is not considered being used for biostratigraphic approaches using neural networks. Not that *P. diamesa* is very similar to its neighboring species like *P. papalis* and *P. phyxis* (Carlsson et al., 2022; Carlsson et al., 2023?). After revising *P. diamesa*, not even one was found. This could be partly because of two reasons; one is that *P. diamesa* is not very abundant within the studied material or it has a very low recall value.

P. phyxis

As *P. sinuosa, P. ampla* and *P. diamesa; P. phyxis* has a low recall value, only resulting in one specimen being identified here. This class also has a low accuracy since species like *P. papalis, P. diamesa* and *P. ampla* are more or less mistaken as *P. phyxis*.

R? auriculaleporis

This class has a high recall, most *R*? *auriculaleporis* is recognized in this class along with other classes, such as *Dictyomitra* Zittel 1876 species. It can be used for semi-automatic identification with some revision by both taxonomists and non-taxonomists.

R? biauritum

This class is rare in our samples, and it looks like *R? auriculaleporis*, just having fewer and more sparsely separated pores. They still have a high recall value and could be used for automatic image classification with some revisions made even for non-taxonomists.

Sethochytris triconiscus Haeckel, 1887

Easily recognizable, having both a high accuracy and recall value. Therefore, it can be used for automatic image recognition by anyone.

T. scolopax

Conversely to *P. lazari*, *T. scolopax* has a little lower recall value since some of its specimens are misclassified as *P. lazari*. However, not a lot of *P. lazari* are recognized as *T scolopax*.

Zealithapium mitra (Ehrenberg, 1875)

Despite using a low or high confidence score, a lot of images are misclassified to this class making it difficult to use for automatic image recognition. Even if this class has a high recall value, one must delete numerous of images that are not *Z. mitra*.

5.4.3 Time and Effort

While training a neural network requires substantial data and significant effort for supervised labeling, this challenge becomes particularly pronounced in the case of radiolarians with close similarities to others, and this dataset, in particular, has too many morphologically similar species.

The process of imaging, on the other hand, is largely automated, requiring minimal manual intervention. However, it is worth noting that accumulating a substantial number of images, around 90,000 images from 23 samples in this study is a bit time-consuming. One sample including 8 coverslips, would take approximately 10 hours, including the time of setting the microscope and saving images. Nevertheless, once you have gathered all the images, the analysis process is fast and consistent if you already have a trained CNN.

5.5 Conclusion

This work concludes that many of these species trained by a CNN can be applied for real biostratigraphic work and interpretations.

Just consider that the sediment slides are prepared in a way to avoid touching each other and forming lumps. Therefore, not a lot of sediment was used. This makes it possible for the CNN

to sometimes not detect rarely occurring species at certain intervals. Some species are better trained than others because they have both a high recall and a high accuracy and could be completely used for automatic identification by non-taxonomists. Other species could be used for automatic recognition with some revisions, while some classes need to be better trained, especially those with a low recall value.

Here, we were sometimes able to expand the biostratigraphic ranges for some classes, while in other cases we were not able to detect certain species intervals due to either poorly trained or the unavailability to capture rare occurring species from preparing the sediment slides.

Acknowledgments

This study has been partly funded by the French government through the program "Investissements d'avenir" (I-ULNE SITE/ANR-16-IDEX-0004 ULNE) managed by the National Research Agency. Funding from the European Union's Horizon 2020 research and innovation program under the Marie Sklodowska-Curie grant agreement no. 847568 was also received and it was also supported by UMR 8198 Evo-Eco-Paléo and IRCICA (CNRS and Univ. Lille USR-3380).

Many thanks to Sylvie Regnier for helping to prepare samples by removal of non-siliceous material.

CHAPTER 6: General Conclusions and Perspectives

6.1 General conclusions

This thesis has aimed to attempt to simplify or automate Eocene radiolarian micropaleontological work by using artificial neural networks for automatic image recognition; is here presented how each chapter answers the following questions, which were presented at the beginning of this thesis:

Q.1. How can neural network learning achieve equal accuracy in the identification of Middle Eocene radiolarian species as an Eocene specialist in radiolarian taxonomy?

Q.2. How can machine learning techniques be compared to classical morphometric analysis?

Q.3. How can Spiking Neural Networks (SNNs) be compared with other deep learning methods?

Q.4 Can a trained neural network on a whole middle Eocene tropical radiolarian assemblage be directly applied for example to biostratigraphic or palaeoceanographic studies?

The second chapter focused on the training of biostratigraphically significant species within the genus *Podocyrtis*. In this study, we exclusively used specimens that belonged to clearly defined specific species, carefully excluding any intermediate or uncertain forms. The training accuracy achieved with MobileNet v.1 reached approximately 91%. We systematically categorized images into three distinct groups: one for unbroken and clear images; another for images with only broken parts removed; and a third group that retained all variations. We tested this approach with newly acquired samples from different locations, yielding accuracy rates ranging from 59% to 77%, depending on the specific dataset used. Nevertheless, it is worth noting that

in the majority of cases, when specimens were not assigned to their true class, specimens were mistakenly classified to at least a closely related species within the same evolutionary lineage.

If a neural network has undergone comprehensive training on morphologically distinct species, there is a strong likelihood that it can achieve accuracy levels comparable to those of an experienced middle Eocene radiolarian taxonomist. One notable advantage of a well-trained neural network is its ability to quickly analyse thousands of images almost instantaneously. Moreover, it can effectively differentiate between closely related and visually similar forms, provided that all such forms have been sufficiently trained. Using a pre-trained CNN for biostratigraphic applications is not only efficient but also enjoyable, as the machine learning model takes charge of specimen identification and counting for each interval. The images just have to be verified how they were categorized into different classes. Nonetheless, it is worth noting that the image acquisition process, even when using an automated microscope, can be time-consuming when dealing with a substantial size of samples.

In a follow-up study conducted by a Master 2 student that I supervised (Francisco Pinto) we dealt with the question of intermediate forms. Our case study focused on intermediate forms between the species *Podocyrtis chalara* and *Podocyrtis goetheana*, representing the end members of the *Lampterium* lineage. We identified two distinct types of intermediate forms. These intermediates, alongside *P. goetheana* and *P. chalara*, underwent attempted identification using various machine learning algorithms based on neural networks, which directly processed image data as input. This approach was compared to Linear Discriminant Analysis (LDA), which utilizes morphometric data for identification purposes. Both LDA and neural networks were able to recognize the intermediate form that is close to *P. goetheana*

(*Podocyrtis* sp. cf. *P. goetheana*). However, both methods were also unable to separate the intermediate form that is nearest to *P. chalara*, *Podocyrtis* sp. cf. *chalara*, since the neural network often misclassified it as *P. chalara*, and the LDA showed that *P. chalara* and *Podocyrtis* sp. cf. *chalara* both share the same morphospace.

By trying to address the second question (Q.2) we achieved similar results when comparing machine learning techniques to classical morphometric analyses, even in the case of distinguishing intermediate forms between two species. Although LDA yielded slightly higher accuracy, the overall outcome remained consistent. The key distinction lies in the methodology. Traditional morphological approaches demand manual measurement of specific values for each individual specimen. In contrast, neural networks simply require the input of images, but to ensure their efficiency, a substantial number of images representing each unique individual, including variations in, for example, rotation is essential. Fortunately, this can be easily obtained using an automated microscope.

To continue on this third chapter, the neural networks used were both traditional CNNs and SNNs. Spiking neural networks are not commonly used for image recognition as widely as CNNs; they are known to perform less well when it comes to accuracy for image recognition, but in recent years they have improved their performance. Apart from being energy efficient and memory saving, SNNs closely imitate natural neural networks and are more suitable to process spatio-temporal data. In this chapter, we could see that using STDP-SNN Network gave a near result of what the trained VGG16 CNN architecture obtained. However, the use of a SuperSpike-SNN Network was less successful, part of it could be explained that they used fewer layers than VGG16, and generally have a smaller network.

To answer the third question (Q.3), there are some SNNs like the STDP-Network, that can give nearly as accurate results as a CNN. The STDP-Network also have a training speed of about 5-6 minutes and VGG16 had in comparison to that a training speed of about 17-18 minutes while VGG16 only performed better by approximately 2 %.

In the fourth chapter, we tackled the comprehensive analysis of radiolarians and other particles, resulting in a substantial dataset containing 12,217 images divided into 96 trained classes, effectively representing a wide range of objects encountered. We achieved an overall training accuracy of approximately 86% and a testing accuracy of about 76%. Within these 96 classes, 39 were categorized at the species level, including biostratigraphically important and frequently encountered species. Our primary objective was to closely examine these species to explore potential applications, which would also provide us with an answer to the last (Q.4) question. Our findings revealed that many of these species could indeed serve in biostratigraphic analyses, where the presence or absence of a class is the key requirement. However, our ultimate goal was to employ these species for relative abundance studies. Unfortunately, we encountered a significant barrier, as a considerable number of specimens appeared in different levels of broken forms and some images consisted of particles that were occasionally lumped together with several specimens, making it difficult to count the exact number of radiolaria. Also, some species or classes exhibited low recall values, indicating a risk of misclassification into multiple other classes. Consequently, these classes were not considered optimal for use in relative abundance studies.

The fifth chapter extends the discussions from the fourth chapter, exploring the practical application of a trained CNN in a biostratigraphic context. We draw parallels with an earlier study conducted on the same Site 1260, essentially building upon prior research. Our findings

in this chapter largely confirmed those of the earlier study, revealing the consistent presence of similar species within specific intervals. However, there were instances where the trained CNN helped extend the biostratigraphic ranges and offered a more comprehensive perspective. Conversely, in cases where certain species were less frequent, the CNN was less effective in detecting their presence within intervals where they should theoretically exist.

6.2 Perspectives

Not only is the focus of the neural networks themselves important but data collection and preparation play a huge part in this, as well. Suppose we start from the beginning with the sample preparation. In that case, it is very important to do this step properly to retrieve the best samples, free from any breakage, or without having unwanted elements like organic or calcareous fossil remains present, and to carefully dissolve and remove clay particles and lumped sediments using an ultrasonic bath.

When preparing slides, striking the right balance is crucial. It is essential to ensure that the sediment is sufficiently distributed to prevent excessive clumping, yet not so sparse that it becomes inefficient when photographing under an automatic microscope also missing rare existing radiolarian species. Currently, we use 1.2 x 1.2 cm slides primarily because the automatic microscope we use generates composite images that are too large for our segmentation software to handle directly and they have to be converted into grayscale to fit the segmentation tool. Producing larger images with lower resolution is not the optimal solution, as it compromises image quality for the sake of accommodating larger images.

A more favourable solution would be the accessibility to more robust softwares capable of handling larger images. Additionally, extending the duration of the automatic microscope's operation to capture a larger area, which doesn't require hourly monitoring, would be a more efficient approach.

Another critical consideration affects the choice of mounting media, as it significantly impacts image contrast, each with its own advantages and weaknesses. When it comes to automatic image recognition, opting for a mounting medium that yields well-contrasted images is ideal. Examples of such media include Canada Balsam (Fig. 6.1A) and Norland epoxy (Fig. 6.1B). In this study, Norland epoxy has been the preferred choice due to its simplicity, rapid drying with a UV lamp, and its ability to produce high-contrast images. However, one limitation of this product is the occasional presence of trapped bubbles and bubble inclusions within radiolarian specimens, which can obviously affect the results of neural networks.

On the other hand, Canada Balsam (Fig. 6.1A) offers the highest image quality, but it has other drawbacks such as slow settlement, an elevated price, and a tendency to become yellow and fade with age (Ravikumar et al., 2019). In contrast, the Eukitt mounting media (Fig. 6.1C) fails to deliver well-contrasted images, although it avoids the issue of bubble inclusions. However, it is not recommended for photographic purposes due to its limitations in contrast enhancement.



Figure 6.1. Radiolarian images A) with Canada Balsam as mounting media, B) with Norland epoxy as monutning media, and C) with Eukitt as mounting media.

In addition to improving neural networks, the creation of robust datasets is equally important. When making these datasets, it is essential to acknowledge that subtle distinctions often differentiate one species from another, and several species may bear an outstanding resemblance. In such cases, it may be necessary to group various similar-looking species, especially when their taxonomic affinity remains uncertain. Conversely, it becomes critically important to collect a substantial number of images for closely resembling species, enabling the neural network to eventually distinguish between these classes.

Furthermore, if neither human experts nor any morphometric approaches can reliably differentiate between two distinct species or face challenges in doing so, it is reasonable to expect the neural network to encounter similar difficulties. Clarity is the key when determining the images assigned to each class; the inclusion of incorrect or uncertain specimens in different classes must be avoided. It's often more practical to initially exclude uncertain specimens, with the possibility of later dividing them into distinct groups as more images are collected and clear distinguishing characteristics emerge.

Additionally, it is preferable to minimize the focus on non-radiolarian elements or other unnecessary particles. Broken fragments, bubbles, and other particles are very common, or even more common than the radiolaria themselves. Overutilizing such data in datasets can lead to overtraining, where the neural network's classification is dominated by these non-relevant features. Hence, careful consideration should be given to have a balance in the dataset composition to ensure effective model training.

One crucial aspect to keep in mind is the incomplete and sometimes unreliable radiolarian taxonomy. Radiolarians can be grouped differently based on their resemblance to taxa that are

not closely related, thereby posing a challenge in their taxonomic classification. Despite extensive efforts to organize radiolarians through morphology and morpho-molecular approaches, the reliability of taxonomy remains questionable, except at the species level. Furthermore, there exists a substantial number of undescribed radiolarian species, and many forms are difficult to describe. These forms often bear strong resemblances to multiple other species and may exist as intermediate forms linking the gap between some established species.

References

- Adobe Systems Incorporated. (2022). Photoshop (Version 23.0) [Computer software]. Retrieved from <u>https://www.adobe.com/products/photoshop.html</u>.
- Aitchison, J. C., Suzuki, N., Caridroit, M., Danelian, T., and Noble, P. (2017). Paleozoic radiolarian biostratigraphy. *In:* Danelian, T., Caridroit, M., Noble, P., and Aitchison, J. C. (Eds.), Catalogue of Paleozoic Radiolarian Genera. Geodiversitas, 39, 503–531. Scientific Publications of the Muséum National d'Histoire Naturelle, Paris. https://doi.org/10.5252/g2017n3a5.

Anderson, O. R. (1983). Radiolaria. Springer Science and Business Media.

- Armstrong, H. A., and Brasier, M. D. (2005). Microfossils. Second Edition. Blackwell Publishing. <u>https://doi.org/10.1002/9780470750995</u>.
- Balouek, D., Carpen Amarie, A., Charrier, G., Desprez, F., Jeannot, E., Jeanvoine, E., Lèbre, A., Margery, D., Niclausse, N., Nussbaum, L., Richard, O., Perez, C., Quesnel, F., Rohr, C., and Sarzyniec, L., (2013). Adding virtualization capabilities to the Grid'5000 testbed: *In:* Cloud Computing and Services Science, Communications in Computer and Information Science). Springer International Publishing, 367, 3-20. https://doi.org/10.1007/978-3-319-04519-1_1.
- Beaufort, L., de Garidel Thoron, T., Mix, A. C., and Pisias, N. G. (2001). ENSO-like forcing on Oceanic Primary Production during the late Pleistocene. Science, 293, 2440–2444. <u>https://doi.org/10.1126/science.293.5539.2440</u>.
- Beaufort, L., and Dollfus, D. (2004). Automatic recognition of coccoliths by dynamical neural networks. Marine Micropaleontology, 51, 57-73. https://doi.org/10.1016/j.marmicro.2003.09.003.

- Boltovskoy, D., Anderson, O. R., and Correa, N. M. (2017). Handbook of the Protists. Springer, Cham, 731–763. https://doi.org/10.1007/978-3-319-28149-0_19.
- Bourel, B., Marchant, R., de Garidel-Thoron, T., Tetard, M., Barboni, D., Gally, Y., and Beaufort, L. (2020). Automated recognition by multiple convolutional neural networks of modern, fossil, intact and damaged pollen grains. Computers and Geosciences, 140, 104498. <u>https://doi.org/10.1016/j.cageo.2020.104498</u>.
- Brocher, J. (2022). biovoxxel/BioVoxxel-Toolbox: BioVoxxel Toolbox v2.5.3. Zenodo. https://doi.org/10.5281/zenodo.5986129.
- Bolli, H. M., Saunders, J. B., and Perch-Nielsen, K. (1985). Eds., Plankton Stratigraphy. Cambridge University Press. ISBN 0 521 23576 6.
- Burki, F., Roger, A. J., Brown, M. W., and Simpson, A. G. B. (2020). The New Tree of Eukaryotes. Trends Ecol Evol, 35 (1), 43-55. https://doi.org/10.1016/j.tree.2019.08.008.
- Campbell, A.S. (1954). Radiolaria. In R.C. Moore (Ed.), Treatise on Invertebrate Paleontology. Lawrence, KS: Geological Society of America, University of Kansas Press, 11-195.
- Carlsson, V. (2022). Podocyrtis Image Dataset [data set]. Recherche Data Gouv. https://doi.org/10.57745/G7CHQL.
- Carlsson, V. (2022). Codes for image preparation and MobileNet CNN [code]. Recherche Data Gouv. https://doi.org/10.57745/J4YL4I.
- Carlsson, V. (2023). Middle Eocene Radiolarian image dataset from ODP Leg 207 (Demerara Rise) [data set]. Recherche Data Gouv. <u>https://doi.org/10.57745/E9YXW6</u>.
- Carlsson, V., Danelian, T., Boulet, P., Devienne, P., Laforge, A., and Renaudie, J. (2022). Artificial intelligence applied to the classification of eight middle Eocene species of the

genus Podocyrtis (polycystine radiolaria). Journal of Micropalaeontology, 41, 165–182. https://doi.org/10.5194/jm-41-165-2022.

- Carlsson, V., Danelian, T., Tetard, M., Meunier, M., Boulet, P., Devienne, P., and Ventalon, S. (2023). Convolutional neural network application on a new middle Eocene radiolarian dataset. Marine Micropaleontology, 183, 102268. https://doi.org/10.1016/j.marmicro.2023.102268.
- Carvalho, L. E., Fauth, G., Baecker Fauth, S., Krahl, G., Moreira, A. C., Fernandes, C. P., and von Wangenheim, A. (2020). Automated Microfossil Identification and Segmentation using a Deep Learning Approach. Marine Micropaleontology, 158, 101890. https://doi.org/10.1016/j.marmicro.2020.101890.
- Cavalier-Smith, T., and Chao, E., and Oates, B. (2004). Molecular phylogeny of Amoebozoa and the evolutionary significance of the unikont Phalansterium. European Journal of Protistology EUR J. PROTISTOL, 40, 21-48. https://doi.org/10.1016/j.ejop.2003.10.001.
- Chollet, F. (2015). Keras [code]. GitHub. URL: <u>https://github.com/fchollet/keras</u> (Last access: 2022-10-21).
- Clark, B. L., Campbell, A. S. (1942). Eocene radiolarian faunas from the Mt Diablo area, California. Geological Society of America, Special Papers, 39, (1) 112. https://doi.org/10.1130/spe39-p1.
- Cortes, C., and Vapnik, V. (1995). Support-vector networks. Machine Learning, 20, 273–297. https://doi.org/10.1007/BF00994018.
- Danelian, T., Asatryan, G., Galoyan, G., Sosson, M., Sahakyan, L., Caridroit, M., Avagyan, A. (2012). Geological history of ophiolites in the Lesser Caucasus and correlation with the

Izmir-Ankara-Erzincan suture zone: insights from radiolarian biochronology. Bulletin de la Société Géologique de France, 183, 331–342. https://doi.org/10.2113/gssgfbull.183.4.331.

- Danelian, T., Baudin, F., Gardin, S., Masure, E., Ricordel, C., Fili, I., Meçaj, T. and Muska, K. (2007). The record of mid Cretaceous oceanic anoxic events from the Ionian zone of southern Albania. Revue de Micropaléontologie, 50, 225–237. https://doi.org/10.1016/j.revmic.2007.06.004.
- Danelian, T., and Johnson, K. G. (2001). Patterns of biotic change in Middle Jurassic to Early Cretaceous Tethyan radiolaria. Marine Micropaleontology, 43, 239-260. https://doi.org/10.1016/S0377-8398(01)00029-9.
- Danelian, T., Le Callonnec, L., Erbacher, J., Mosher, D. C., Malone, M. J., Berti, D., Bice, K.
 L., Bostock, H., Brumsack, H. -J., Forster, A., Heidersdorf, F., Henderiks, J., Janecek,
 T. J., Junium, C., MacLeod, K., Meyers, P. A., Mutterlose, J. H., Nishi, H., Norris, R.
 D., Ogg, J. G., and Glatz, C. (2005). Preliminary results on Cretaceous-Tertiary tropical
 Atlantic pelagic sedimentation (Demerara Rise, ODP Leg 207). Comptes Rendus
 Geoscience, 337, 609–616. https://doi.org/10.1016/j.crte.2005.01.011.
- Danelian, T., and Macleod, N. (2019). Morphometric Analysis of Two Eocene Related Radiolarian Species of the Podocyrtis (Lampterium) Lineage. Paleontological Research, 23, 314–330. https://doi.org/10.2517/2019PR007.
- Danelian, T., and Monnet, C. (2021). Early Paleozoic radiolarian plankton diversity and the Great Ordovician Biodiversification Event. Earth-Science Reviews, 218, no. 103672. <u>https://doi.org/10.1016/j.earscirev.2021.103672</u>.

- Danelian, T., Saint Martin, S., and Blanc-Valleron, M. -M. (2007). Middle Eocene radiolarian and diatom accumulation in the equatorial Atlantic (Demerara Rise, ODP Leg 207):
 Possible links with climatic and palaeoceanographic changes. C. R. Palevol., 6, 103–114. <u>https://doi.org/10.1016/j.crpv.2006.08.002</u>.
- Danelian, T., Zambetakis-Lekkas, A., Galoyan, G., Sosson, M., Asatryan, G., Hubert, B., and Reconstructing Grigoryan, A. (2014).Upper Cretaceous (Cenomanian) paleoenvironments in Armenia based on Radiolaria and benthic Foraminifera; implications for the geodynamic evolution of the Tethyan realm in the Lesser Caucasus. Palaeogeography, Palaeoclimatology, Palaeoecology, 413, 123–132. https://doi.org/10.1016/j.palaeo.2014.03.011.
- Davies, M., Wild, A., Orchard, G., Sandamirskaya, Y., Fonseca Guerra, G. A., Joshi, P., Plank,
 P., and Risbud, S. (2021). Advancing Neuromorphic Computing with Loihi: A Survey of Results and Outlook, Proc. IEEE, 109, 911–934, https://doi.org/10.1109/JPROC.2021.3067593.
- De Lima, R. P., Welch, K. F., Barrick, J. E., Marfurt, K. J., Burkhalter, R., Cassel, M., and Soreghan, G. S. (2020). Convolutional Neural Networks As An Aid To Biostratigraphy And Micropaleontology: A Test On Late Paleozoic Microfossils. Palaios, 35, 391–402. <u>https://doi.org/10.2110/palo.2019.102</u>.
- De Wever, P., Dumitrica, P., Caulet, J. P., Nigrini, C., and Caridroit, M. (2002). Radiolarians in the Sedimentary Record. London: CRC Press. https://doi.org/10.1201/9781482283181.

- Dollfus, D., and Beaufort, L. (1996). Automatic pattern recognition of calcareous nanoplankton. Proceedings of the Conference on Neural Networks and their Applications (NEURAP 96), 306–311. Marseille, France.
- Dollfus, D., and Beaufort, L. (1999). Fat neural network for recognition of position-normalized objects. Neural Networks, 12, 553-560. <u>https://doi.org/10.1016/S0893-6080(99)00011-</u>
 - <u>8</u>.
- Ehrenberg, C. G. (1846). Uber eine halibiolithische, von Herrn R. Schomburgk entdeckte, vorherrschend aus mikroskopischen Polycystinen gebildete, Gebirgsmasse von Barbados, Bericht uber die zu Bekanntmachung geeigneten Verhandlungen der Koniglichen Preussische Akademie derWissenschaften zu Berlin, Jahre 1846, 382–385.
- Ehrenberg, C. G. (1847). Über die mikroskopischen kieselschaligen Polycystinen als mächtige Gebirgsmasse von Barbados und über das Verhältnis der aus mehr als 300 neuen Arten bestehenden ganz eigenthümlichen Formengruppe jener Felsmasse zu den jetzt lebenden Thieren und zur Kreidebildung, Eine neue Anregung zur Erforschung des Erdlebens. K. Preuss. Akad. Wiss., Berlin, Jahre 1847, 40–60.
- Ehrenberg, C. G. (1854). Mikrogeologie: das Erden und Felsen schaffende Wirken des unsichtbar kleinen selbstständigen Lebens auf der Erde. Leopold Voss, Leipzig, Germany, 374 pp. <u>https://doi.org/10.5962/bhl.title.118752</u>.
- Ehrenberg, C. G. (1874). Grössere Felsproben des Polycystinen-Mergels von Barbados mit weiteren Erläuterungen. K. Preuss. Akad. Wiss., Berlin, Monatsberichte, Jahre 1873, 213–263.
- Ehrenberg, C. G. (1875). Fortsetzung der mikrogeologischen Studien Abhandlungen der königlichen Akademie der Wissenschaften.

- Elbez, H. (2023). Code for: Morphometrics and machine learning discrimination of the middle Eocene radiolarian species Podocyrtis chalara, P. goetheana and their morphological intermediates. Deposited at <u>https://archive.softwareheritage.org/browse/directory/cc7d8ef1505299a208adcde597a</u> <u>98d90b0ca47d6/</u> (Last access 2023-08-24).
- Erbacher, J., Mosher, D. C., Malone, M. J., Berti, D., Bice, K. L., Bostock, H., Brumsack, H.-J., Danelian, T., Forster, A., Glatz, C., Heidersdorf, F., Henderiks, J., Janecek, T. R., Junium, C., Le Callonnec, L., MacLeod, K., Meyers, P. A., Mutterlose, H. J., Nishi, H., Norris, R. D., Ogg, J. G., O'Regan, A. M., Rea, B., Sexton, P., Sturt, H., Suganuma, Y., Thurow, J. W., Wilson, P. A. and Wise Jr., S. W. (2004). Proceedings of the Ocean Drilling Program, Initial Reports, 207. College Station, TX: Ocean Drilling Program. https://doi.org/10.2973/odp.proc.ir.207.2004.
- Falez, P. (2019). Improving Spiking Neural Networks Trained with Spike Timing DependentPlasticity for Image Recognition [PhD thesis]. Université de Lille.
- Gonçalves, A. B., Souza, J. S., Silva, G. G., Cereda, M. P., Pott, A., Naka, M. H., and Pistori, H. (2016). Feature Extraction and Machine Learning for the Classification of Brazilian Savannah Pollen Grains. PLoS One, 11, e0157044. https://doi.org/10.1371/journal.pone.0157044.
- Haeckel, E. (1887). Report on the Radiolaria collected by H.M.S. Challenger during the years 1873-1876. Report on the Scientific Results of the Voyage of H.M.S. Challenger during the years 1873-1876. 18: 1-1803.
- Hay, W. W., De Conto, R., Wold, C. N., Wilson, K. M., Voigt, S., Schulz, M., Wold-Rossby, A., Dullo, W. -C., Ronov, A. B., Balukhovsky, A. N., and Soeding, E. (1999).

Alternative global Cretaceous paleogeography. *In*: Barrera, E., and Johnson, C. (Eds.), The Evolution of Cretaceous Ocean/Climate Systems, Geological Society of America Special Paper 332, 1-47. <u>https://doi.org/10.1130/0-8137-2332-9.1</u>.

- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep Residual Learning for Image Recognition.
 In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 770-778. <u>https://doi.org/10.48550/arXiv.1512.03385</u>.
- Hijazi, S., Kumar, R., and Rowen, C. (2015). Using Convolutional Neural Networks for Image Recognition. Cadence, Cadence Design Systems Inc., 1–12.
- Hollis, C. J., Pascher, K. M., Sanfilippo, A., Nishimura, A., Kamikuri, S.-I., and Shepherd, C.
 L. (2020). An Austral radiolarian biozonation for the Paleogene. Stratigraphy, 17, 213– 278. <u>https://doi.org/10.29041/strat.17.4.213-278</u>.
- Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., and Adam, H. (2017). Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861. https://doi.org/10.48550/arXiv.1704.04861.
- Hsiang, A.Y., Brombacher, A., Rillo, M.C., Mleneck-Vautravers, M.J., Conn, S., Lordsmith, S., Jentzen, A., Henehan, M. J., Metcalfe, B., Fenton, I. S., Wade, B. S., Fox, L., Meilland, J., Davis, C. V., Baranowski, U., Groeneveld, J., Edgar, K. M., Movellan, A., Aze, T., Dowsett, H. J., Giles Miller, C., Rios, N., and Hull, P. M. (2019). Endless Forams: >34,000 modern planktonic foraminiferal images for taxonomic training and automated species recognition using convolutional neural networks. Paleoceanography and Paleoclimatology, 34, 1157-1177. <u>https://doi.org/10.1029/2019PA003612</u>.

- Itaki, T., Taira, Y., Kuwamori, N., Saito, H., Ikehara, M., and Hoshino, T. (2020). Innovative microfossil (radiolarian) analysis using a system for automated image collection and AI-based classification of species. Scientific Reports. <u>https://doi.org/10.1038/s41598-020-77812-6</u>.
- Kamikuri, S. -I., Moore, T. C., Ogane, K., Suzuki, N., Pälike, H., and Nishi, H. (2012). Early Eocene to early Miocene radiolarian biostratigraphy for the low-latitude Pacific Ocean. Stratigraphy, 9, 77–108. <u>https://doi.org/10.29041/strat.09.1.03</u>.
- Kim, K. (1992). Paleogene radiolarian biostratigraphy from high latitude. South Atlantic Journal of the Paleontological Society of Korea, 8 (1), 24-51.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). ImageNet Classification with Deep Convolutional Neural Networks. Commun. ACM, 60, 84–90. https://doi.org/10.1145/3065386.
- Lampitt, R. S., Salter, I., and Johns, D. (2009). Radiolaria: Major exporters of organic carbon to the deep ocean. Global Biogeochemical Cycles, 23, GB1010. https://doi.org/10.1029/2008GB003221.
- Lazarus, D. (2005). A brief review of radiolarian research. Paläontologische Zeitschrift, 79(1), 183-200. Stuttgart, 31. 3. 2005. <u>https://doi.org/10.1007/BF03021761</u>.
- Lever, J., Krzywinski, M., and Altman, N. (2017). Principal component analysis. Nature Methods, 14, 641–642. <u>https://doi.org/10.1038/nmeth.4346</u>.
- Lloyd, S. P. (1957). Least squares quantization in PCM. Technical Report RR-5497, Bell Lab, September 1957.

- Maass, W. (1997). Networks of spiking neurons: The third generation of neural network models. Neural Networks, 10, 1659–1671. <u>https://doi.org/10.1016/S0893-6080(97)00011-7</u>.
- MacQueen, J. B. (1967). "Some Methods for Classification and Analysis of Multivariate Observations". Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability. 1: 281–297.
- Marchant, R., Tetard, M., Pratiwi, A., Adebayo, M., and de Garidel-Thoron, T. (2020). Automated analysis of foraminifera fossil records by image classification using a convolutional neural network. Journal of Micropalaeontology, 39, 183–202. https://doi.org/10.5194/jm-39-183-2020.
- Masquelier, T., and Thorpe, S. J. (2007). Unsupervised learning of visual features through spike timing dependent plasticity. PLoS computational biology, 3(2), e31. <u>https://doi.org/10.1371/journal.pcbi.0030031</u>.
- Matsuoka, A (2007). Living radiolarian feeding mechanisms: new light on past marine ecosystems. Swiss Journal of Geosciences. 100 (2): 273–279. https://doi.org/10.1007/s00015-007-1228-y.
- Matsuzaki, K. M., Itaki, T., Tada, R., and Kamikuri, S. (2018). Paleoceanographic history of the Japan Sea over the last 9.5 million years inferred from radiolarian assemblages (IODP Expedition 346 Sites U1425 and U1430). Progress in Earth and Planetary Science, 5(1), 54. <u>https://doi.org/10.1186/s40645-018-0204-7</u>.
- Meunier, M., and Danelian, T. (2022). Astronomical calibration of late middle Eocene radiolarian bioevents from ODP Site 1260 (equatorial Atlantic, Leg 207) and refinement

of the global tropical radiolarian biozonation. Journal of Micropalaeontology, 41, 1–27. https://doi.org/10.5194/jm-41-1-2022.

- Meunier, M., and Danelian, T. (2023). Progress in understanding middle Eocene nassellarian (Radiolaria, Polycystinea) diversity; new insights from the western equatorial Atlantic Ocean. Journal of Paleontology, 97, 1–25. <u>https://doi.org/10.1017/jpa.2022.82</u>.
- Mitra, R., Marchitto, T. M., Ge, Q., Zhong, B., Kanakiya, B., Cook, M. S., Fehrenbacher, J. S., Ortiz, J. D., Tripati, A., and Lobaton, E. (2019). Automated species-level identification of planktic foraminifera using convolutional neural networks, with comparison to human performance. Marine Micropaleontology, 148, 1-14. https://doi.org/10.1016/j.marmicro.2019.01.005.
- Moore, T. C. (1971). Radiolaria. *In*: Initial Reports of the Deep Sea Drilling Project, Volume 8, edited by Tracey Jr, J. I., Sutton, G. H., Nesteroff, W. D., Galehouse, J., Von der Borch, C.C., Moore, T., Lipps, J., Haq, U. Z. B. U., and Beckmann, J.P. U.S. Govt. Print. Office, Washington, DC, USA, 727-775. https://doi.org/10.2973/dsdp.proc.8.112.1971.
- Moore, T. C. (1972). Mid-Tertiary Evolution of the Radiolarian Genus Calocycletta. Micropaleontology, 18, 144–152. <u>https://doi.org/10.2307/1484991</u>.
- Nigrini, C.A., Sanfilippo, A., and Moore, T.J., (2005). Radiolarian biostratigraphy and chronology of radiolarian events of ODP Leg 199 sites and EW9709 sediment cores:
 Supplement to: Nigrini, CA et al. (2005): Cenozoic Radiolarian Biostratigraphy: A Magnetobiostratigraphic Chronology of Cenozoic Sequences from ODP Sites 1218, 1219, and 1220, Equatorial Pacific. *In*: Wilson, PA; Lyle, M; Firth, JV (Eds.)

Proceedings of the Ocean Drilling Program, Scientific Results, College Station, TX (Ocean Drilling Program), 199, 1-76. https://doi.org/10.2973/odp.proc.sr.199.225.2005.

- Nigrini, C. (1974). Cenozoic Radiolaria from the Arabian Sea, DSDP Leg 23. *In*: Initial Reports of the Deep Sea Drilling Project, 23, edited by: Whitmarsh, R. B., Weser, O. E., Ali, S., Boudreaux, J. E., Fleisher, R. L., Jipa, D., Kidd, R. B., Mallik, T. K., Matter, A., Nigrini, C., Siddiquie, H. N., and Stoffers, P., U.S. Govt. Print. Office, Washington, DC, USA, 1051–112. <u>https://doi.org/10.2973/dsdp.proc.26.233.1974</u>.
- Obut, O., and Iwata, K. (2000). Lower Cambrian Radiolaria from the Gorny Altai (southern West Siberia). Journal of Geology and Geophysics, 41.
- O'Dogherty, L., Caulet, J., Dumitrica, P., and Suzuki, N. (2021). Catalogue of Cenozoic radiolarian genera (Class Polycystinea). Geodiversitas, 43, 709–1185. https://doi.org/10.5252/geodiversitas2021v43a21.
- Oksanen, J., Blanchet, F.G., Kindt, R., Legendre, P., Minchin, P. R., O'Hara, R. B., Simpson, G.L., Solymos, P., Stevenes, M.H.H., Wagner, H. (2013). Package 'vegan': Community Ecology Package. Software <u>http://CRAN.R-project.org/package=vegan</u>.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
 Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D.,
 Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine Learning in
 Python. Journal of Machine Learning Research, 12, 2825–2830.
- Pehle, C. -G., and Pedersen, J. E. (2021). Norse A deep learning library for spiking neural networks (0.0.5). Zenodo. <u>https://doi.org/10.5281/zenodo.4422025</u>.

- Pinto, F., Carlsson, V., Meunier, M., Van Bocxlaer, B., Elbez, H., Cueille, M., Boulet, P., and Danelian, T. (2023). Data for the: Morphometrics and machine learning discrimination of the middle Eocene radiolarian species Podocyrtis chalara, P. goetheana and their morphological intermediates. Recherche Data Gouv. <u>https://doi.org/10.57745/8KBOFP</u>.
- Pouille, L., Obut, O., Danelian, T., and Sennikov, N. (2011). Lower Cambrian (Botomian) polycystine Radiolaria from the Altai Mountains (southern Siberia, Russia). Comptes Rendus Palevol, 10, 627–633.
- R Core Team. (2022). R: a language and environment for statistical computing: R Foundation for Statistical Computing. <u>https://www.R-project.org</u>.
- Ravikumar, S. Shamala., Surekha, R. and Thavarajah, Rooban. (2014). Mounting media: An overview. Journal of Dr. NTR University of Health Sciences. 3, 1.
 <u>https://doi.org/10.4103/2277-8632.128479</u>.
- Renaudie, J., Danelian, T., Saint Martin, S., Le Callonnec, L., and Tribovillard, N. (2010). Siliceous phytoplankton response to a Middle Eocene warming event recorded in the tropical Atlantic (Demerara Rise, ODP Site 1260A). Palaeogeography, Palaeoclimatology, Palaeoecology, 286, 121–134. <u>https://doi.org/10.1016/j.palaeo.2009.12.004</u>.
- Renaudie, J., Gray, R., and Lazarus, D. (2018). Accuracy of a neural net classification of closely-related species of microfossils from a sparse dataset of unedited images. PeerJ Preprints, 6:e27328v1.

- Renaudie, J. and Lazarus, D. (2013). On the accuracy of paleodiversity reconstructions: a case study in Antarctic Neogene radiolarians. Paleobiology, 39 (3), 491-509. <u>https://doi.org/10.1666/12016</u>.
- Riedel, W. R. (1971). Cenozoic Radiolaria from the western tropical Pacific, Leg 7. Init. Repts. DSDP, 7, 1592–1627.
- Riedel, W. R. and Sanfilippo, A. (1970). Radiolaria, Leg 4, Deep Sea Drilling Project. *In*: Initial Reports of the Deep Sea Drilling Project, 4. <u>https://doi.org/10.2973/dsdp.proc.4.124.1970</u>.
- Riedel, W. R. and Sanfilippo, A. (1978). Stratigraphy and Evolution of Tropical Cenozoic Radiolarians. Micropaleontology, 24, 61–96. <u>https://doi.org/10.2307/1485420</u>.
- Ripley, B., Venables, B., Bates, D. M., Hornik, K., Gebhardt, A., Firth, D. and Ripley, M. B. (2013). Package 'mass': Cran r, 538, 113–120.
- Rohlf, F., and Bookstein, F. (1991). Size and Shape (Book Reviews: Proceedings of the Michigan Morphometrics Workshop). Science, 253, 345-362. https://doi.org/10.1126/science.253.5024.1152.
- Rueckauer, B., Bybee, C., Goettsche, R., Singh, Y., Mishra, J., and Wild, A. (2021). NxTF: An API and Compiler for Deep Spiking Neural Networks on Intel Loihi. arXiv [preprint]. https://doi.org/10.48550/arXiv.2101.04261.
- Sandin, M. (2019). Diversity and Evolution of Nassellaria and Spumellaria (Radiolaria). Protistology. Sorbonne Université. English. NNT : 2019SORUS549ff. tel-03137926

- Sandin, M., Pillet, L., Biard, T., Poirier, C., Bigeard, E., Romac, S., Suzuki, N., and Not, F. (2019). Time Calibrated Morpho-molecular Classification of Nassellaria (Radiolaria).
 Protist, 170 (2), 187-208. <u>http://doi.org/10.1016/j.protis.2019.02.002</u>.
- Sanfilippo, A., and Nigrini, C. (1998). Code numbers for Cenozoic low latitude radiolarian biostratigraphic zones and GPTS conversion tables: Marine Micropaleontology, 33, 109–156. <u>https://doi.org/10.1016/S0377-8398(97)00030-3</u>.
- Sanfilippo, A., and Riedel, W. R. (1970). Post-Eocene "Closed" Theoperid Radiolarians: Micropaleontology, 16, 446–462. <u>http://doi.org/10.2973/gsmicropal.16.4.446</u>.
- Sanfilippo, A., and Riedel, W. R. (1973). Cenozoic Radiolaria (Exclusive of Theoperids, Artostrobiids, and Amphipyndacids) from the Gulf of Mexico, Deep Sea Drilling Project Leg 10. <u>http://doi.org/10.2973/DSDP.PROC.10.119.1973</u>.
- Sanfilippo, A., and Riedel, W. R. (1990). Morphometric Analysis of Evolving Eocene Podocyrtis (Radiolaria) Morphotypes Using Shape Coordinates. *In*: Proceedings of the Michigan Morphometrics Workshop, 345–362.
- Sanfilippo, A., and Riedel, W. R. (1992). The Origin and Evolution of Pterocorythidae (Radiolaria): A Cenozoic Phylogenetic Study: Micropaleontology, 38, 1. https://doi.org/10.2307/1485841.
- Sanfilippo, A., Westberg-Smith, M. J., and Riedel, W. R. (1985). Cenozoic Radiolaria. *In*: Plankton Stratigraphy, 631–712. <u>https://doi.org/10.1017/S2475263000001422</u>.
- Schneider, C. A., Rasband, W. S., and Eliceiri, K. W., (2012). NIH Image to ImageJ: 25 years of image analysis: Nature Methods, 9, 671–675. <u>https://doi.org/10.1038/nmeth.2089</u>.
- Shipboard Scientific Party: Site 612 (1987). *In*: Initial reports of the Deep Sea Drilling Project covering Leg 95 of the cruises of the drilling vessel Glomar Challenger, 95, 313–337,
edited by: Poag, C. W., Watts, A. B., Cousin, M., Goldberg, D., Hart, M. B., Miller, K. G., Mountain, G. S., Nakamura, Y., Palmer, A. A., Schiffelbein, P. A., Schreiber, B. C., Tarafa, M. E., Thein, J. E., Valentine, P. C., Wilkens, R. H., and Turner, K. L., St. John's, Newfoundland, to Ft. Lauderdale, Florida, August–September 1983, Texas A and M University, Ocean Drilling Program, College Station, TX, United States, http://doi.org/10.2973/dsdp.proc.95.103.1987.

- Shipboard Scientific Party: Site 711 (1988). *In*: Proceedings of the Ocean Drilling Program, Initial Reports, Vol. 115, edited by: Backman, J., Duncan, R. A., Peterson, R. C., Baker, A. B., Baxter, A. L., Boersma, A., Cullen, A., Droxler, A. W., Fisk., M. R., Greenough, J.D., Hargraves, R. B., Hempel., P., Hobart, M. A., Hurley, M. T., Johnson., D. A., Macdonald, A. H., Mikkelsen, N., Okada, H., Rio, D., Robinson, S. G., Schneider, D., Swart, P. K., Tatsumi, Y., Vandamme, D., Vilks, G., Vincent, E (Participating Scientists), Peterson, L. C. (Shipboard Staff Scientist), and Barbu, E. M., Ocean Drilling Program, http://doi.org/10.2973/odp.proc.ir.115.110.1988.
- Shipboard Scientific Party: Site 1051 (1998). *In*: Proceedings of the Ocean Drilling Program,
 Initial Reports, 171B, 171–239, edited by: Norris, R. D., Kroon, D., Klaus, A.,
 Alexander, I. T., Bardot, L. P., Barker, C. E., Bellier, J-P., Blome, C. D., Clarke, L. J.,
 Erbacher, J., Faul, K. L., Holmes, M. A., Huber, B. T., Katz, M. E., MacLeod, K. G.,
 Marca, S., Martinez-Ruiz, F. C., Mita, I., Nakai, M., Ogg, J. G., O'Regan, A. M., Rea,
 B., Sexton, P., Sturt, H., Suganuma, Y., Thurow, J. W., Wilson, P. A., Baez, L. A., and
 Kapitan-White, E., College Station, TX,
 http://doi.org/10.2973/odp.proc.ir.171b.105.1998.

- Shipboard Scientific Party: Site 1260 (2004). *In*: Proceedings of the Ocean Drilling Program, Initial Reports, 207, edited by: Erbacher, J., Mosher, D. C., Malone, M. J., Berti, D., Bice, K. L., Bostock, H., Brumsack, H.-J., Danelian, T., Forster, A., Glatz, C., Heidersdorf,F., Henderiks, T., Janecek, T. R., Junium, C., Le Callonnec, L., MacLeod, K., Meyers, P. A., Mutterlose, H. J., Nishi, H., Norris, R. D., Ogg, J. G., O'Regan, A. M., Rea, B., Sexton, P., Sturt, H., Suganuma, Y., Thurow, J. W., Wilson, P. A., and Wise Jr., S. W., Ocean Drill. Program, College Station, TX, USA, http://doi.org/10.2973/odp.proc.ir.207.107.2004.
- Simonyan, K., and Zisserman, A. (2015). Very Deep Convolutional Networks for Large-Scale Image Recognition: ICLR, San Diego, California, USA, <u>https://arxiv.org/abs/1409.1556</u>.
- Souza, A. L., Eilert, V. M. P., Souza Lima Fidalgo, T., dos Reis, I. P., Vilela, C. G., and Mendonça Filho, J. G. (2022). Late middle Eocene to early Oligocene radiolarian bios.
 J. Micropalaeontology, 41, 1–27. <u>https://doi.org/10.1016/j.marmicro.2021.102038</u>.
- Souza, A. L., Eilert, V. M. P., Souza Lima Fidalgo, T., and Mendonça Filho, J. G. (2017). Early to middle Eocene radiolarian biostratigraphy for the mid-latitude South Atlantic Ocean, Site 356, DSDP LEG 39, Mar. Micropaleontol., 136, 66–89. http://doi.org/10.1016/j.marmicro.2017.09.004.
- Suzuki, N., and Yoshiaki, A. (2011). Radiolaria: Achievements and unresolved issues: Taxonomy and cytology. Plankton and Benthos Research. 6. 69-91. https://doi.org/10.3800/pbr.6.69.
- Suzuki, N., O'Dogherty, L., Caulet, J. -P., and Dumitrica, P. (2021). A new integrated morphoand molecular systematic classification of Cenozoic radiolarians (Class Polycystinea) –

suprageneric taxonomy and logical nomenclatorial acts. Geodiversitas, 43(15), 405-573. <u>http://doi.org/10.5252/geodiversitas2021v43a15</u>.

- Takahashi, K. (1987). Radiolarian flux and seasonality: climatic and El Nino response in the subarctic Pacific, 1982–1984. Global Biogeochem Cycles 1, 213–231. https://doi.org/10.1029/GB001i003p00213.
- Tavanaei, A., Ghodrati, M., Kheradpisheh, S. R., Masquelier, T., and Maida, A., (2019). Deep learning in spiking neural networks. Neural Networks, v. 111, p. 47-63. https://doi.org/10.1016/j.neunet.2018.12.002.
- Tetard, M., Marchant, R., Cortese, G., Gally, Y., de Garidel-Thoron, T., and Beaufort, L., (2020). Technical note: A new automated radiolarian image acquisition, stacking, processing, segmentation and identification workflow. Climate of the Past, v. 16, p. 2415–2429. https://doi.org/10.5194/cp-16-2415-2020.
- Tetard, M., Monnet, C., Noble, P., and Danelian T., (2017). Biodiversity patterns of Silurian Radiolaria. Earth-Science Reviews, v. 173, p. 77-83. <u>https://doi.org/10.1016/j.earscirev.2017.07.011</u>.
- Trubovitz, S., Lazarus, D., Renaudie, J., Noble, P. (2020). Marine plankton show threshold extinction response to Neogene climate change. Nature Communications, 11, Article no. 5069. <u>https://doi.org/10.1038/s41467-020-18879-7</u>.
- Van der Walt, S., Schönberger, J. L., Nunez-Iglesias, J., Boulogne, F., Warner, J. D., Yager, N., Gouillart, E., Yu, T., and the scikitimage contributors. (2014). scikit-image: Image processing in Python. PeerJ, 2, 453. <u>http://doi.org/10.7717/peerj.453</u>.
- Van Rossum, G. and Drake, F. L. (2009). Python 3 Reference Manual, Scotts Valley, CA, CreateSpace, ISBN: 978-1-4414-1269-0.

- Venables, W. N., and Ripley, B. D., (2002). Random and Mixed Effects, in Venables, W.N. and Ripley, B.D., eds., Modern Applied Statistics with S, Statistics and Computing: New York, NY, Springer, 271–300. <u>https://doi.org/10.1007/978-0-387-21706-2_10</u>.
- Vrielynck, B., Bonneau, M., Danelian, T., Cadet, J. P., Poisson, A. (2003). New insights on the Antalya Nappes in the apex of the Isparta angle: The Isparta Cay unit revisited. Geol. J. 38, 283-293. <u>https://doi.org/10.1002/gj.956</u>.
- Watanabe, M., Kawagata, S., Aita, Y., Suzuki, N., and Kamikuri, S. (2022). Changes in morphological parameters of the radiolarian Lampterium lineage from the middle Eocene in the tropical Pacific. Marine Micropaleontology, v. 173, 102125. https://doi.org/10.1016/j.marmicro.2022.102125.
- Westerhold, T., and Röhl, U. (2013). Orbital pacing of Eocene climate during the Middle Eocene Climate Optimum and the chron C19r event: Missing link found in the tropical western Atlantic. Geochem. Geophy. Geosy., 14, 4811–4825. http://doi.org/10.1002/ggge.20293.
- Yang, Z. R., and Yang, Z., (2014). Artificial Neural Networks. *In*: Comprehensive Biomedical Physics: Elsevier, p. 1–17.
- Zenke, F., and Ganguli, S., (2018). SuperSpike: Supervised Learning in Multilayer Spiking Neural Networks. Neural Computation, v. 30, p. 1514–1541. http://doi.org/10.1162/neco a 01086.
- Zittel, K. A. (1876b). Uber einige fossile Radiolarien aus der norddeutschen Kreide. Zeitschrift der Deutschen Geologischen Gesellschaft. 28: 75-86.

Appendices

Tropical Atlantic middle Eocene Radiolaria catalogue from the Demerara Rise, ODP Leg 207

Radiolarians here are used to train CNNs for automatic image classification and have been balanced to both fit an acceptable taxonomic framework as well as being adjusted to be trained with a high accuracy.

List of all classes

ACANTHODESMIOIDEA ACANTHODESMIIDAE Eucoronis hertwigi group ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE cephalospyridid group A ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE cephalospyridid group B ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE cephalospyridid group C ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE cephalospyridid group D ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Dendrospyris stylophora ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Elaphospyris didiceros group ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Unknown cephalospyridid group ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Liriospyris clathrata group ACANTHODESMIOIDEA CEPHALOSPYRIDIDAE Petalospyris anthocyrtoides group ACANTHODESMIOIDEA juvenil or broken Acanthodesmioidea ACANTHODESMIOIDEA Smooth surfaced Acanthodesmioidea ACANTHODESMIOIDEA Spinose Acanthodesmioidea ACANTHODESMIOIDEA STEPHANIIDAE Zygocircus spp AMPHIPYNDACOIDEA AMPHIPYNDACIDAE Amphipternis cf. clava ARCHAEODICTYOMITROIDEA ARCHAEODICTYOMITRIDAE Dictyomitra parva ARTOSTROBIOIDEA ARTOSTROBIIDAE artostrobiid group ARTOSTROBIOIDEA ARTOSTROBIIDAE Artostrobus quadriporus ARTOSTROBIOIDEA ARTOSTROBIIDAE Dictyoprora amphora group ARTOSTROBIOIDEA ARTOSTROBIIDAE Dictyoprora mongolfieri ARTOSTROBIOIDEA ARTOSTROBIIDAE Dictyoprora spp ARTOSTROBIOIDEA ARTOSTROBIIDAE Siphocampe spp ARTOSTROBIOIDEA ARTOSTROBIIDAE Siphocample pupa group ARTOSTROBIOIDEA ARTOSTROBIIDAE Tricolapsa spp ARTOSTROBIOIDEA RHOPALOSYRINGIIDAE Rhopalosyringium auriculaleporis ARTOSTROBIOIDEA RHOPALOSYRINGIIDAE Rhopalosyringium biauritum Background blurry Broken Broken diatoms CARPOCANIOIDEA CARPOCANIIDAE Carpocanopsis ornata group CYCLADOPHOROIDEA CYCLADOPHORIDAE Cycladophora spatiosa group Diatoms Double Edges EUCYRTIDIOIDEA EUCYRTIDIIDAE Eucyrtidium spp EUCYRTIDIOIDEA EUCYRTIDIIDAE Stichopterygium microporum EUCYRTIDIOIDEA EUCYRTIDIIDAE Stichopterygium picus Foraminifera HELIOSATURNALOIDEA AXOPRUNIDAE Axoprunum sp A LITHELIOIDEA SPONGURIDAE Middourium group LITHOCHYTRIDOIDEA BEKOMIDAE Dictyophimus craticula LITHOCHYTRIDOIDEA LITHOCHYTRIDIDAE lithochytridid group LITHOCHYTRIDOIDEA LITHOCHYTRIDIDAE Lithochytris vespertilio LITHOCHYTRIDOIDEA LITHOCHYTRIDIDAE Lychnocanium babylonis group LITHOCHYTRIDOIDEA LITHOCHYTRIDIDAE Lychnocanoma bajunensis LITHOCHYTRIDOIDEA LITHOCHYTRIDIDAE Sethochytris triconiscus LITHOCYCLIOIDEA LITHOCYCLIIDAE Lithocyclia ocellus LITHOCYCLIOIDEA PHACODISCIDAE Periphaena decora NASSELLARIA group A NASSELLARIA group B NASSELLARIA group C Particles PLAGIACANTHOIDEA CERATOCYRTIDAE Ceratocyrtis spp PLAGIACANTHOIDEA DIMELISSIDAE Lithomelissa ehrenbergii group PLAGIACANTHOIDEA DIMELISSIDAE Lithomelissa macroptera PLAGIACANTHOIDEA group A PLAGIACANTHOIDEA group B PLAGIACANTHOIDEA group C PLAGIACANTHOIDEA LOPHOPHAENIDAE Lophophaena radians group PLAGIACANTHOIDEA LOPHOPHANIDAE cf Lophophaena simplex PLAGIACANTHOIDEA LOPHOPHANIDAE Lophophaena capito group PLAGIACANTHOIDEA LOPHOPHANIDAE Lophophaena spp

PLAGIACANTHOIDEA PLAGIACANTHIDAE Rhabdolithis pipa PLECTOPYRAMIOIDEA Incertae familiae Zealithapium mitra PLECTOPYRAMIOIDEA Incertae familiae Zealithapium spp PLECTOPYRAMIOIDEA PLECTOPYRAMIDIDAE plectopyramidid group PTEROCORYTHOIDEA LOPHOCYRTIIDAE Aphetocyrtis zamenhofi PTEROCORYTHOIDEA LOPHOCYRTIIDAE Apoplanius spp PTEROCORYTHOIDEA LOPHOCYRTIIDAE lophocyrtiid group A PTEROCORYTHOIDEA LOPHOCYRTIIDAE lophocyrtiid group B PTEROCORYTHOIDEA LOPHOCYRTIIDAE Lophocyrtis alauda PTEROCORYTHOIDEA LOPHOCYRTIIDAE Lophocyrtis barbadense PTEROCORYTHOIDEA PTEROCORYTHIDAE Albatrossidium spp PTEROCORYTHOIDEA PTEROCORYTHIDAE Calocyclas hispida PTEROCORYTHOIDEA PTEROCORYTHIDAE Calocycloma ampulla PTEROCORYTHOIDEA PTEROCORYTHIDAE Phormocyrtis embolum PTEROCORYTHOIDEA PTEROCORYTHIDAE Phormocyrtis lazari PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Lampterium chalara PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Lampterium goetheana PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Lampterium mitra group PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Lampterium puellasinensis PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Lampterium sinuosa PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Podocyrtis papalis PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Podocyrtoges ampla PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Podocyrtoges diamesa PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Podocyrtoges phyxis PTEROCORYTHOIDEA PTEROCORYTHIDAE Podocyrtis Podocyrtopsis apeza PTEROCORYTHOIDEA PTEROCORYTHIDAE pterocorythid group PTEROCORYTHOIDEA PTEROCORYTHIDAE Theocyrtis scolopax PTEROCORYTHOIDEA PTEROCORYTHIDAE Theocyrtis spp PTEROCORYTHOIDEA THEOCOTYLIDAE Pterocodon apis PTEROCORYTHOIDEA THEOCOTYLIDAE Theocorys spp PTEROCORYTHOIDEA THEOCOTYLIDAE Theocotyle spp PTEROCORYTHOIDEA THEOCOTYLIDAE Thyrsocyrtis Pentalocorys krooni PTEROCORYTHOIDEA THEOCOTYLIDAE Thyrsocyrtis Pentalocorys triacantha PTEROCORYTHOIDEA THEOCOTYLIDAE Thyrsocyrtis Thyrsocyrtis rhizodon PTEROCORYTHOIDEA THEOPERIDAE Eusyringium fistuligerum group PTEROCORYTHOIDEA THEOPERIDAE Eusyringium lagena PTEROCORYTHOIDEA THEOPERIDAE Rhopalocanium ornatum group PTEROCORYTHOIDEA THEOPERIDAE Rhopalocanium sphinx PYLOBOTRYDOIDEA PYLOBOTRYDIDAE pylobotrydid group Spicules SPUMELLARIA Discoidal biconvex spumellarians with equatorial spines SPUMELLARIA Discoidal spongy spumellarians SPUMELLARIA Discoidal spongy spumellarians with rod shaped spines SPUMELLARIA Discoidal spumellarians with irregular concentric rings SPUMELLARIA Discoidal spumellarians with regular concentric rings SPUMELLARIA Multi armed spongy spumellarians SPUMELLARIA Small discoidal spumellarians with concentric rings SPUMELLARIA Spherical spumellarians with large pores SPUMELLARIA Spherical spumellarians with radial spines SPUMELLARIA Two armed spongy spumellarians STICHOPILIOIDEA STICHOPILIIDAE Lophoconus antilope STYLOSPHAEROIDEA STYLOPHAERIDAE Spongatractus pachystylus STYLOSPHAEROIDEA STYLOSPHAERIDAE Stylosphaera coronata group THEOPILIOIDEA ANTHOCYRTIDIDAE Anthocyrtis mespilus group

Nassellarians

Superfamily **ARCHAEODICTYOMITROIDEA** Pessagno, 1976 Family **Archaeodictyomitridae** Pessagno, 1976

Genus Dictyomitra Zittel, 1876



Dictyomitra parva (Kim, 1992)

Superfamily AMPHIPYNDACOIDEA Riedel, 1967 Family Amphipyndacidae Riedel, 1967 Genus Amphipternis Foreman, 1973



Amphipternis cf. clava (Ehrenberg, 1874)

This species differs from *Amphipternis clava* (Ehrenberg) in having more than two postabdominal segments, numerous closely spaced pores arranged in transverse row, and in being more inflated distally.

Not trained by the CNN

50 µm

Dictyomitra parva (Kim)

Eucyrtidium parva Kim, 1992, p. 43, pl. 2, figs. 7, 8.

Theoperid gen. et sp. indet. Johnson, 1974, pl. 3, fig. 12.

Siphocampe elizabethae (Clark and Campbell): Nigrini, 1977, p. 256, pl. 3, fig. 6.

Archaeodictyomitra ? sp. Takemura, 1992, p. 744, pl. 3, figs. 1, 2.

Dictyomitra amygdala Shilov, 1995, p. 126, pl. 1, fig. 4-6b.

Dictyoprora ? amygdala (Shilov): Suzuki et al., 2009, p. 263, pl. 18, fig. 3.

Siphocampe ? amygdala (Shilov): Hollis et al., 2020, pl. 8, figs. 11, 12; Meunier and Danelian, 2022, p. 22, pl. 1, fig. 1.

Amphipternis cf. clava Ehrenberg, 1874

cf. *Lithocampe* ? *clava* Ehrenberg, 1874, p. 238; 1876, p. 76, pl. 4, fig. 3; Ogane et al., 2009, pl. 22, figs. 2a–2c.

Amphipternis clava (Ehrenberg): Foreman, 1973, p. 430, pl. 7, figs. 16, 17; pl. 9, fig. 2.

Superfamily EUCYRTIDIOIDEA Ehrenberg, 1846 Family Eucyrtidiidae Ehrenberg, 1846 Genus Eucyrtidium Ehrenberg, 1847



Eucyrtidium spp.

We included here *Eucyrtidium levisaltarix* (Meunier and Danelian, 2023), *Eucyrtidium montiparum* Ehrenberg, 1874 and other *Eucyrtidium* species with more than two post-abdominal segments, absence of feet or lateral horns and a cylindrical test outline.

Genus Stichopterygium Haeckel, 1882



Stichopterygium ? microporum (Ehrenberg, 1874)



Stichopterygium ? picus (Ehrenberg, 1874) *Not trained by the CNN*

Stichopterygium ? microporum (Ehrenberg)

Eucyrtidium microporum Ehrenberg, 1874, p. 230; 1876, p. 72, pl. 11, fig. 20; Ogane et al., 2009, pl. 6, figs. 5a–5c, pl. 85, figs. 5a–5f.

Stichopodium ? microporum (Ehrenberg): Petrushevskaya and Kozlova, 1972, p. 548, pl. 25, figs. 4–6; Funakawa et al., 2006, p. 37, pl. P13, figs. 3a–4b.

Stichopterygium ? picus (Ehrenberg)

Eucyrtidium picus Ehrenberg, 1874, p. 232; 1876, p. 72, pl. 11, fig. 1; Ogane et al., 2009, pl. 86, figs. 6a–6g.

Artostrobiid gen. et sp. indet. Riedel and Sanfilippo, 1977, pl. 9, fig. 16.

Superfamily **PLECTOPYRAMIDOIDEA** Haecker, 1908 Family **Plectopyramididae** Haecker, 1908



Plectopyramidid group

Elongated conical two segmented plectopyramidids that are wide open distally and display well-arranged pore rows on the test. In some species the pores are almost square, in others they are oval, or circular to subcircular in outline. Many species have a very tiny spherical cephalis that appears poreless and bears sometimes a horn that may be tiny or large tubular or thick conical.

Family *Incertae familiae* Genus *Zealithapium* O'Connor, 1999



Zealithapium mitra (Ehrenberg, 1874)



Zealithapium spp.

Included in this class are *Zealithapium plegmacantha* (Riedel and Sanfilippo, 1970) and *Zealithapium anoectum* (Riedel and Sanfilippo, 1970).

Zealithapium mitra (Ehrenberg)

Cornutella mitra Ehrenberg, 1874, p. 221; 1876, p. 68, pl. 2, figs. 8. *Lithapium ? mitra* (Ehrenberg): Riedel and Sanfilippo, 1970, p. 520, pl. 4, figs. 6, 7. *Zealithapium mitra* (Ehrenberg): O'Connor, 1999, p. 5.

Superfamily CARPOCANIOIDEA Haeckel, 1882 Family Carpocaniidae Haeckel, 1882 Genus Carpocanopsis Riedel and Sanfilippo, 1971







50 µm

Carpocanopsis ornata (Ehrenberg, 1874) group

Superfamily ARTOSTROBIOIDEA Riedel, 1967 Family Artostrobiidae Riedel, 1967 Genus Artostrobus Petrushevskaya, 1967



50 μm Artostrobus quadriporus Bjørklund, 1976 Not trained by the CNN

Carpocanopsis ornata (Ehrenberg) group

Cryptoprora ornata Ehrenberg, 1874, p. 222; 1876, p. 68, pl. 5, fig. 8; Ogane et al., 2009, pl. 6, figs. 2a–2c, pl. 83, figs. 5a–6d.

Cryptocarpium ornatum (Ehrenberg): Sanfilippo and Riedel, 1992, p. 6, 36, pl. 2, figs. 18–20.

Remarks: The great morphological disparity observed in *Carpocanopsis ornata* (Ehrenberg, 1874) leads us to consider it as a group of closely related species.

Artostrobus quadriporus Bjørklund

Artostrobus quadriporus Bjørklund, 1976, p. 1125, pl. 23, figs. 15-21; Hull, 1996, p. 137, pl. 4, fig. 12.

non ? Artostrobus quadriporus Lazarus and Pallant, 1989, p. 362, pl. 5, figs. 16, 17.

Superfamily **ARTOSTROBIOIDEA** Riedel, 1967 Family **Artostrobiidae** Riedel, 1967 Genus *Dictyoprora* Haeckel, 1881







Dictyoprora mongolfieri (Ehrenberg, 1854)









Dictyoprora amphora (Haeckel, 1887) group

Dictyoprora with several tightly packed pores. In some specimens the test surface is rather rough and the pore size increases at the distal end of the thorax. Pores are quincuncially arranged, not strictly arranged in transverse and longitudinal rows like in *D. mongolfieri* (Ehrenberg, 1854).



Dictyoprora spp.

All *Dictyoprora* specimens that were uncertain to belong to either *D. mongolfieri* or *D. amphora* group were placed here. In this group there are also many other species such as *Dictyoprora ovata* (Haeckel, 1887), *Dictyoprora pirum* (Ehrenberg, 1874), *Dictyoprora curta* (Clark and Campbell, 1942), *Dictyoprora urceolus* Haeckel, 1887 and many othergunknown forms.

Dictyoprora mongolfieri (Ehrenberg)

Eucyrtidium mongolfieri Ehrenberg, 1854, pl. 36, fig. 18; 1874, p. 230; 1876, p. 72, pl. 10, fig. 3.

Dictyoprora mongolfieri (Ehrenberg): Nigrini, 1977, p. 250, pl. 4, fig. 7; Funakawa et al., 2006, p. 17, pl. P2, figs. 5a-6b.

Dictyoprora amphora (Haeckel) group

Dictyocephalus amphora Haeckel, 1887, p. 1305, pl. 62, fig. 4.

Theocampe amphora (Haeckel) group: Foreman, 1973, p. 431, pl. 8, figs. 7, 9–13.

Dictyoprora amphora (Haeckel) group: Nigrini, 1977, p. 250, pl. 4, figs. 1, 2; Funakawa et al., 2006, p. 16, 17, pl. P2, figs. 1a–2b.

Genus Siphocampe Haeckel, 1882



Siphocampe pupa (Ehrenberg, 1874) group

Two segmented artostrobiids with a large thorax and small horizontally arranged tightly packed pores.



Siphocampe pollen





Siphocampe pachyderma

50 µm Siphocampe acephala

Siphocampe spp.

We included in this group three species, *Siphocampe pollen* Meunier and Danelian, 2023; *Siphocampe pachyderma* (Ehrenberg, 1874), and *Siphocampe acephala* (Ehrenberg, 1874).

Siphocampe pupa (Ehrenberg, 1861) group

Eucyrtidium pupa Ehrenberg 1861, p. 768; 1873a, p. 311; 1873b, pl. 7, fig. 16; Suzuki et al., 2009, pl. 55, figs. 8a–8c.

Remarks: The great morphological disparity observed in *Siphocampe pupa* (Ehrenberg, 1874) leads us to consider it as a group of closely related species.

Genus Tricolocapsa Haeckel, 1882



"Tricolocapsa" spp.

Two segmented form often bearing a short conical horn. The cephalis is much smaller than the thorax, which seems inflated. Some specimens bear feet; it is doubtful whether all of the specimens here actually belong to the family Artostrobiidae.





Artostrobiid group

A lump group of artostrobiids with several segments and a curvy outline belonging to the genera *Botryostrobus* and *Siphocampe*. Not trained by the CNN

Family **Rhopalosyringiidae** Empson-Morin, 1981 Genus *Rhopalosiringium* Campbell and Clark, 1944



Rhopalosyringium ? auriculaleporis (Clark and Campbell, 1942)



Rhopalosyringium ? biauritum (Ehrenberg, 1874)

Rhopalosyringium ? auriculaleporis (Clark and Campbell)

Lophophaena (Lophophaenula) auriculaleporis Clark and Campbell, 1942: p. 76; pl. 8, figs. 20, 27–29; Blueford, 1988, p. 246, pl. 3, figs. 1–3.

Artobotrys auriculaleporis (Clark and Campbell): Petrushevskaya and Kozlova, 1979, p. 137, fig. 515; Meunier and Danelian, 2022, p. 18, pl. 1, fig. 7.

Rhopalosyringium ? biauritum (Ehrenberg)

Eucyrtidium biauritum Ehrenberg, 1874, p. 226; 1876, p. 70, pl. 10, figs. 7, 8; Ogane et al., 2009, pl. 18, figs. 8a–d, pl. 20, figs. 1a–2b, 6.

Eucyrtidium bicorne Ehrenberg, 1874, p. 226; 1876, p. 70, pl. 11, fig. 7; Ogane et al., 2009, pl. 20, figs. 3a, 3b, 4a–4c, 5a–5c.

Lophocyrtis biaurita (Ehrenberg): Haeckel, 1887, p. 1411; Riedel and Sanfilippo, 1978, p. 70, pl. 6, fig. 13.

Artobotrys biaurita (Ehrenberg): Petrushevskaya and Kozlova, 1979, p. 136; Meunier and Danelian, 2022, p. 18, pl. 1, fig. 6.

Superfamily ACANTHODESMIOIDEA Haeckel, 1862 Family Acanthodesmiidae Haeckel, 1882

Genus Eucoronis Haeckel, 1882



Eucoronis hertwigi (Bütschli, 1882) group Not trained by the CNN

Eucoronis hertwigi (Bütschli) group

Acanthodesmia hertwigi Bütschli, 1882, pl. 32, fig. 9a-9c.

Eucoronis hertwigii (Bütschli) group: Petrushevskaya and Kozlova 1972, p. 533, pl. 41, figs. 15-17; Kamikuri, 2015, pl. 13, fig. 18.

Family Cephalospyrididae Haeckel, 1882

Genus Dendrospyris Haeckel, 1882



Dendrospyris stylophora (Ehrenberg, 1874)

Genus Elaphospyris Haeckel, 1882



Elaphospyris ? didiceros (Ehrenberg, 1874) group

We place here all cephalospyridids with one apical horn, two antapical horns diverging laterally at ca. 45° and variable number of feet.

Dendrospyris stylophora (Ehrenberg)

Ceratospyris stylophora Ehrenberg, 1874, p. 220; 1876, p. 66, pl. 20, fig. 10; Ogane et al., 2009, pl. 38, fig. 6a–6c, pl. 39, figs. 6a, 6b.

Dendrospyris stylophora (Ehrenberg): Goll, 1968, p. 1423, pl. 173, figs. 21-24, text-fig. 8.

Elaphospyris ? didiceros (Ehrenberg) group

Ceratospyris didiceros Ehrenberg, 1874, p. 218; 1876, p. 66, pl. 21, fig. 6; Ogane et al., 2009, pl. 39, figs. 1a–1c.

Giraffospyris didiceros (Ehrenberg): Goll, 1969, p. 332, pl. 60, figs. 5-7, 9, text-fig. 2.

Giraffospyris didiceros (Ehrenberg) group: Riedel and Sanfilippo, 1970, pl. 5, figs. 3-5.

Dendrospyris didiceros (Ehrenberg) group: Petrushevskaya and Kozlova, 1972, p. 532, pl. 40, fig. 12.

Genus Liriospyris Haeckel, 1882



Liriospyris clathrata (Ehrenberg, 1854) group

We include in this group cephalospyridids with a large variation in terms of test thickness, size of pores, size and number of feet and spines.

Genus Petalospyris Ehrenberg, 1846



"Petalospyris anthocyrtoides" Bütschli, 1882 group

Two-segmented nassellarians with a weak collar stricture line between the cephalis and thorax. Cephalis appears to be poreless. One row of pores is developed under the lumbar stricture from which extend several spiny feet or spines.

Liriospyris clathrata (Ehrenberg) group

Dictyospyris clathrus Ehrenberg, 1854, pl. 36, fig. 25B

Dictyospyris clathrata Ehrenberg, 1874, p. 224; 1876, p. 68, pl. 19, fig. 7; Bütschli, pl. 32, figs. 10a, b.

Liriospyris clathrata (Ehrenberg): Goll, 1968, p. 175, figs. 12, 13, 16, 17; Ling, 1975, p. 726, pl. 7, figs. 6–9.

Remarks: There is large morphological variation in this group which leads us to determine this as a group of species.

"Petalospyris anthocyrtoides" (Bütschli, 1882) group

Petalospyris anthocyrtoides Bütschli 1882, p.532, 538, pl.32, figs.19a-c.

Remarks: There are some morphological variations within this group which leads us to determine this as a group of species.



Unknown cephalospyridid group

A rather sphaerical morphotype ca. 100 μ m in diameter, with a thick walled test and rough test surface. Pores are variable in size, but mostly large. The cephalic lobes are bulbous and not always identical to each another. Feet or teeth-like structures taper distally.



Cephalospyridid group A

Cephalospyridids with or without horn and variable number of feet. Most cephalospyrids included in this class are unknown or rare species, but some of them may be identified as *Dorcadospyris ombros* Nigrini et al., 2006 and *Triceraspyris articulata* (Ehrenberg, 1874).



Cephalospyridid group B

We placed here all thick-walled *Desmospyris* species with a partly developed thorax. The encountered specimens vary in shape. Some of them may be identified as *Desmospyris obtusus* (Bütschli, 1882) group. Many specimens display cephalic lobes of irregular shape; some specimens bear constricted small feet or teeth-like structures and others have tiny apical horns. The pore size and arrangement is often irregular. We also included here some forms belonging to the genus *Petalospyris*, alike *Petalospyris confluens* Ehrenberg, 1874.







Cephalospyridid group C

Cephalospyridids with smooth test, bearing a short conical horn and 2 - 3 straight feet placed with a ca. 45° angle with respect to the base of the cephalis.



Cephalospyridid group D

Included species in this class are *Petalospyris argiscus* Ehrenberg, 1874 and *Thamnospyris fragoides* (Sanfilippo and Riedel, 1973).

Family **Stephaniidae** Haeckel, 1882 Genus *Zygocircus* Bütschli, 1882



Zygocircus spp.

We placed here all *Zygocircus* species found in our material, including some that resemble to *Z. butschlii*, but they are smaller and lack a small corner ring.



Smooth surfaced Acanthodesmioidea

Bilobate forms with a D-shaped sagittal ring, with or without feet and some of which may display a tiny horn. Some have a rough test surface, while others are rather smooth. Encountered specimens in our material belong to either Acanthodesmiidae or Cephalospyridae.



Ceratospyris fibula

Spinose Acanthodesmioidea

Bilobate forms that display well developed spines, horns and feet, some of which may be identified at the species level (i.e. *Ceratospyris fibula* Ehrenberg, 1874).







50 µm
Superfamily **THEOPILIOIDEA** Haeckel, 1882 Family **Theopilioidea** Haeckel, 1882 Genus Anthocyrtis Ehrenberg, 1846



Anthocyrtis mespilus (Ehrenberg, 1847) group

Two-segmented nassellarians with a weak collar stricture line between the cephalis and thorax. Cephalis seems to be mostly poreless. Several spiny feet or spines replacing the abdomen with maximum one pore row under the lumbar stricture.

Superfamily STICHOPILOIDEA Haeckel, 1882 Family Stichopiliidae Haeckel, 1882 Genus Lophoconus Haeckel, 1887



Lophoconus antilope (Ehrenberg, 1874) *Not trained by the CNN*

Anthocyrtis mespilus Ehrenberg group

Anthocyrtis mespilus Ehrenberg, 1847, p. 55, fig. 9; 1854, pl. 36, fig. 13; 1876, p. 66, pl. 6, figs. 4, 5; Ogane et al., 2009, pl. 50, figs. 2a, 2b, pl. 80, figs. 5a–5e, pl. 81, figs. 1a–2d; Kamikuri, 2015, pl. 10, figs. 4a–6.

Anthocyrtis furcata Ehrenberg, 1874, p. 216; 1876, p. 64, pl. 6, fig. 2; Funakawa et al., 2006, p. 38, pl. P13, figs. 5a, 5b; Ogane et al., 2009, pl. 80, figs. 4a–4f.

Remarks: There are a lot of morphological variabilities within this group which leads us to determine this as a group of species.

Lophoconus antilope (Ehrenberg, 1874)

Lophoconus antilope (Ehrenberg, 1874), p. 308.

Superfamily **PLAGIACANTHOIDEA** Hertwig, 1879 Family **Ceratocyrtidae** Petrushevskaya, 1981

Genus Ceratocyrtis Bütschli, 1882



Ceratocyrtis spp.

We placed here various *Ceratocyrtis* species, Including *Ceratocyrtis ampliata* (Ehrenberg, 1874).

Family **Dimelissidae** Petrushevskaya, 1981 Genus *Lithomelissa* Ehrenberg, 1847





Lithomelissa ehrenbergi Buetschli, 1882 group

We placed here all two-segmented lophophaeniids with two to three long bladed lateral wings extending from the thorax. A bladed horn is present in most specimens. Pores are small, numerous and arranged irregularly.





Lithomelissa macroptera Ehrenberg, 1874

Lithomelissa ehrenbergi Bütschli, 1882 group

Lithomelissa ehrenbergi Bütschli, 1882, p. 517, pl. 33, fig. 21a, 21b.

Remarks: There are a lot of morphological variabilities within this group which leads us to determine this as a group of species.

Lithomelissa macroptera Ehrenberg, 1874

Lithomelissa macroptera Ehrenberg, 1874, p. 241; 1876, p. 78, pl. 3, figs. 8-10; Ogane et al., 2009, pl. 4, figs. 3a–3c, 4a–4d, 7a–7d, pl. 19, figs. 6a–6d; Kamikuri, 2015, pl. 11, figs. 16, 17.

Family **Lophophaenidae** Haeckel, 1882 Genus *Lophophaena* Ehrenberg, 1847





Lophophaena capito (Ehrenberg, 1874) group

Lophophaeniids with a distinctly inflated large cephalis, which displays a constricted bottleneck at its base. A small conical apical horn is often present. *Not trained in the CNN*



Lophophaena radians Ehrenberg, 1874 group

We placed here all lophophaenids with several apical horns or spines that are sometimes distally and medially attached like a trabecule. There are mainly three different morphotypes present here, all of them considered as part of the species variability, including typical *L. radians* with large pores and a rather distinct collar stricture. There are morphotypes that have less distinct collar strictures and pores of equal size that are smaller and more sparsely distributed.

Lophophaena capito Ehrenberg group

Lophophaena capito Ehrenberg, 1874, p. 242; 1876, p. 78, pl. 8, fig. 6; Ogane et al., 2009, pl. 19, figs. 8a–8c, pl. 34, figs. 3a–3c, pl. 79, figs. 2a–2c.

Lophophaena ? capito Ehrenberg group: Petrushevskaya and Kozlova, 1972, p. 535, pl. 33, figs. 20-23.

Lophophaena capito Ehrenberg group: Funakawa et al., 2006, p. 20, Pl. P3, figs. 3-4.

Lophophaena radians Ehrenberg group

Lophophaena radians Ehrenberg, 1874, p. 243; 1876, p. 78, pl. 8, figs. 7–9; Ogane et al., 2009, pl. 3, figs. 3a–3e, 5a–5d, pl. 79, figs. 4a–4c; Funakawa et al., 2006, p. 20, pl. P3, figs. 5a–6b; Kamikuri, 2015, pl. 13, figs. 30a, 30b; Meunier and Danelian, 2022, p. 21, pl. 3, fig. 1.



Lophophaena **spp.** Two segmented lophophaenids bearing a short horn.

Family Phaenocalpididae Haeckel, 1887

Genus Rhabdolithis Ehrenberg, 1847



Rhabdolithis pipa Ehrenberg, 1854

Rhabdolithis pipa Ehrenberg

Rhabdolithis pipa Ehrenberg, 1854, pl. 36, fig. 59; 1876, p. 159, pl. 1, fig. 27; Sanfilippo and Riedel, 1973, p. 529, pl. 18, figs. 12–16, pl. 33, figs. 9, 10.



Plagiacanthoidea group A Unknown plagicanthoidea.



Plagiacanthoidea group B Two-segmented forms with two to three long feet and a strong bladed horn..



Plagiacanthoidea group C

Two-segmented forms that are often without any spines and only exceptionally display several small feet.

Superfamily **PYLOBOTRYDOIDEA** Haeckel, 1882 Family **Pylobotrydidae** Haeckel, 1882



Pylobotrydid group

Two-segmented nassellarians with a cephalis that is subdivided into ante-, eu-,and postcephalic lobes.

Superfamily CYCLADOPHOROIDEA Suzuki, 2019 Family Cycladophoridae Suzuki, 2019 Genus Cycladophora Ehrenberg, 1846



Cycladophora spatiosa (Ehrenberg, 1874) group

Three segmented cycladophorids that are overall conical in shape and display a globose cephalis, a campanulate thorax and a truncated conical abdomen that is wide open distally. Pores are rather large and well-organised in rows.

Cycladophora spatiosa Ehrenberg group

Cycladophora spatiosa Ehrenberg, 1847, p. 48; 1874, p. 222; 1876, p. 68, pl. 18, figs. 5, 6; Haeckel, 1887, p. 1379; Ogane et al., 2009, pl. 9, figs. 6a, 6b, pl. 87, figs. 4a, 4b; Kamikuri, 2015, pl. 10, figs. 7a, 7b.

Cycladophora spatiosa (Ehrenberg) group: Funakawa et al., 2006, p. 38, pl. P13, figs. 7a-8b.

Anthocyrtella spatiosa (Ehrenberg): Petrushevskaya and Kozlova, 1972, p. 541, pl. 33, figs. 1–3.

Superfamily **LITHOCHYTRIDOIDEA** Ehrenberg, 1846 Family **Bekomidae** De Wever et al., 2001

Genus Dictyophimus Ehrenberg, 1847



Dictyophimus craticula Ehrenberg, 1874 *Not trained by the CNN*





Lithochytris vespertilio Ehrenberg, 1874

Genus Sethochytris Haeckel, 1882



Sethochytris triconiscus Haeckel, 1887

Dictyophimus craticula Ehrenberg

Dictyophimus craticula Ehrenberg, 1874, p. 223; 1876, p. 68, pl. 5, figs. 4, 5; Sanfilippo and Riedel, 1973, p. 529, pl. 19, fig. 1 (*partim.*); Ogane et al., 2009, pl. 21, fig. 5, pl. 36, figs. 1a–1f, pl. 37, figs. 2a–4; Meunier and Danelian, 2022, p. 19, pl. 3, fig. 6.

Lithochytris vespertilio Ehrenberg

Lithochytris vespertilio Ehrenberg, 1874, p. 239; 1876, p. 76, pl. 4, fig. 10; Riedel and Sanfilippo, 1970, p. 518, pl. 9, figs. 8, 9; Ogane et al., 2009, pl. 45, figs. 1a–3e; Meunier and Danelian, 2022, p. 20, pl. 3, fig. 4.

Sethochytris triconiscus Haeckel

Sethochytris triconiscus Haeckel, 1887, p. 1239, pl. 57, fig. 13; Riedel and Sanfilippo, 1970, p. 528, pl. 9, figs. 6; Meunier and Danelian, 2022, p. 22, pl. 1, fig. 15.

Genus Lychnocanoma Haeckel, 1887



Genus Lychnocanium Ehrenberg, 1846



Lychnocanium babylonis (Clark and Campbell, 1942) **group** We placed here all lithochytridoiids with a conical horn that covers most of or the entire cephalis and bear three smooth non-bladed feet.



Lithochytridid group

We included in this class all two-segmented lithochytridids with a horn, other than *Lychnocanissa bajunensis*, *Lychnocanium babylonis*, or *Lithochytris* spp. bearing porous feet. An example of an identified species is *Lychnocanium carinatum* Ehrenberg, 1874. 232

Lychnocanoma bajunensis Renz

Lychnocanoma bajunensis Renz 1984, p. 459, pl. 1, figs. 4-6.

Lychnocanium babylonis (Clark and Campbell) group

? Lychnocanium tribulus Ehrenberg, 1874, p. 245; 1876, p. 80, pl. 7, fig. 1.

Dictyophimus (Dictyophimium) babylonis Clark and Campbell, 1942, p. 67, pl. 9, figs. 32, 36.

Sethochytris babylonis (Clark and Campbell) group: Riedel and Sanfilippo, 1970, p. 528, pl. 9, figs. 1–3.

Superfamily **PTEROCORYTHOIDEA** Haeckel, 1882 Family **Lophocyrtiidae** Sanfilippo and Caulet in De Wever et al., 2001 Genus *Aphetocyrtis* Sanfilippo and Caulet, 1998



Aphetocyrtis zamenhofi Meunier and Danelian, 2023

Genus Apoplanius Sanfilippo & Caulet 1998



Apoplanius spp.

Included forms are *Apoplanius keraspera* (Sanfilippo and Caulet, 1998) and *Apoplanius aspera* (Ehrenberg, 1874).

Aphetocyrtis zamenhofi Meunier and Danelian

Theocorys sp. Riedel and Sanfilippo, 1977, pl. 7, fig. 9. *Aphetocyrtis zamenhofi* Meunier and Danelian, 2023, p. 17, figs. 6.1–6.4

Genus Lophocyrtis Haeckel, 1887





Lophocyrtis alauda (Ehrenberg, 1874)



Lophocyrtis barbadense (Ehrenberg, 1874) *Not trained by the CNN*

Lophocyrtis alauda (Ehrenberg)

Eucyrtidium alauda Ehrenberg, 1874, p. 225; 1876, p. 70, pl. 9, fig. 4; Ogane et al., 2009, pl. 49, figs. 1a–1e.

Lophocyrtis (Lophocyrtis ?) cf.. semipolita: Sanfilippo and Caulet, 1998, p. 10, pl. 4, fig. 7.

Lophocyrtis (Lophocyrtis ?) alauda (Ehrenberg): Meunier and Danelian, 2022, p. 20, pl. 1, fig. 16.

Lophocyrtis barbadense (Ehrenberg)

Pterocanium barbadense Ehrenberg, 1874, p. 254; 1876, pl. 17, fig. 6.

Lophocyrtis (Lophocyrtis ?) barbadense (Ehrenberg): Sanfilippo and Caulet, 1998, p. 8, pl. 4, figs. 9, 10a, 10b; Funakawa et al., 2006, p. 26, pl. P8, figs. 4a-5b.



Lophocyrtiid group A

Rare lophocyrtiids. At least one in this group is *Apoplanius klydus* (Sanfilippo and Caulet, 1998)

Not trained by the CNN



Lophocyrtiid group B

The lophocyrtiids in this group have fewer and larger pores, some forms have a very long and curvy abdomen. Included forms are *Aphetocyrtis? columboi* Meunier and Danelian, 2023 and *Lophocyrtis attenuata* (Ehrenberg, 1874). It is possible that there are other species included in this group as well.

Family **Pterocorythidae** Haeckel, 1882 Genus *Albatrossidium* Sanfilippo and Riedel, 1992



Albatrossidium spp.

Pterocorythids with a relatively large cephalis. Pores of the thorax similar to those of the abdomen. Abdomen in this class is often short or absent. A well-defined collar constriction exists between thorax and cephalis. Presence of a broad-based apical horn.



Calocycloma ampulla (Ehrenberg, 1854) *Not trained by the CNN*

Calocyclas hispida (Ehrenberg)

Anthocyrtis hispida Ehrenberg, 1874, p. 216; 1876, p. 64, pl. 8, fig. 2; Ogane et al., 2009, pl. 2, figs. 7a–9c, pl. 50, figs. 4a, 4b.

Cycladophora hispida (Ehrenberg): Riedel and Sanfilippo, 1970, p. 529, pl. 10, fig. 9; Moore, 1971, p. 741, pl. 4, figs. 6, 7.

Calocyclas hispida (Ehrenberg): Foreman, 1973, p. 434, pl. 1, figs. 12–15, pl. 9, fig. 18; Sanfilippo and Blome, 2001, p. 210, fig. 6g.

Calocycloma ampulla (Ehrenberg)

Eucyrtidium ampulla Ehrenberg, 1854, p. 21, pl. 36, figs. 15a–15c; 1874, p. 225; 1876, p. 70, pl. 10, figs. 11, 12; Ogane et al., 2009, pl. 60, figs. 1a–3d.

Calocycloma ? ampulla (Ehrenberg): Riedel and Sanfilippo, 1970, p. 524, pl. 6, fig. 1.

Calocycloma ampulla (Ehrenberg): Foreman, 1973, p. 434, pl. 1, figs. 1–5, pl. 9, fig. 20.

Genus Phormocyrtis Haeckel, 1887



Phormocyrtis embolum Ehrenberg, 1874



Phormocyrtis lazari Meunier and Danelian, 2023



50 µm

Phormocyrtis embolum (Ehrenberg)

Eucyrtidium embolum Ehrenberg, 1874, p. 228; 1876, p. 70, pl. 10, fig. 5; Johnson, 1974, p. 548, pl. 4, fig. 5 (*partim.*); Nigrini, 1974, p. 1068, pl. 1H, figs. 4, 5; Ogane et al., 2009, pl. 22, figs. 6a–6c.

Phormocyrtis embolum (Ehrenberg): Haeckel, 1887, p. 1369; Riedel, 1957, p. 88, pl. 3, fig. 6 (*partim.*); Meunier and Danelian, 2022, p. 21, pl. 2, fig. 12.

Phormocyrtis lazari Meunier and Danelian

Phormocyrtis lazari Meunier and Danelian, 2023.

Genus *Podocyrtis* Ehrenberg, 1847 Subgenus *Lampterium* Sanfilippo and Riedel, 1992



Podocyrtis (Lampterium) chalara Riedel and Sanfilippo, 1970



Podocyrtis (Lampterium) goetheana (Haeckel, 1887)

Podocyrtis (Lampterium) chalara Riedel and Sanfilippo

Podocyrtis (Lampterium) chalara Riedel and Sanfilippo, 1970, p. 535, pl. 12, figs. 2, 3; Riedel and Sanfilippo, 1978, p. 71, pl. 8, fig. 3; Sanfilippo et al., 1985, p. 697, fig. 30.11

Podocyrtis (Lampterium) goetheana (Haeckel)

Cycladophora goetheana Haeckel, 1887, p. 1376, pl. 65, fig. 5.

Podocyrtis (Lampterium) goetheana (Haeckel): Riedel and Sanfilippo, 1970, p. 535; Sanfilippo et al., 1985 p. 697, fig. 30.12; Nigrini et al., 2005, p. 45, pl. P5, figs. 11, 12.



Podocyrtis (Lampterium) mitra Ehrenberg, 1854 group

Included here are also specimens belonging to *Podocyrtis (Lampterium) trachodes* Riedel and Sanfilippo, 1970.



"Podocyrtis (Lampterium) puellasinensis" Ehrenberg, 1874

Two segmented pterocorytoiids (cephalis and thorax with a missing abdomen) that are conical in shape. It's possible that some of these are or parts of *Podocyrtis Lampterium* spp.



Podocyrtis (Lampterium) sinuosa Ehrenberg, 1874

Podocyrtis (Lampterium) mitra Ehrenberg group

Podocyrtis mitra Ehrenberg, 1854, p. 21, pl. 36, fig. 20.

Podocyrtis (Lampterium) mitra Ehrenberg: Riedel and Sanfilippo, 1970, p. 534, pl. 11, figs. 5, 6; Riedel and Sanfilippo, 1978, pl. 8, fig. 7; Sanfilippo et al., 1985, p. 698, fig. 30.10; Sanfilippo and Blome, 2001, p. 215, figs. 10a, 10b.

Podocyrtis (Lampterium) trachodes Riedel and Sanfilippo, 1970, p. 535, pl. 11, fig. 7; pl. 12, fig. 1; Sanfilippo et al., 1985, p. 699, fig. 30.14; Sanfilippo and Blome; 2001, p. 215, fig. 10c.

"Podocyrtis (Lampterium) puellasinensis" Ehrenberg

Podocyrtis puella sinensis Ehrenberg, 1874, p. 252; 1876, p. 82, pl. 14, fig. 3.

Podocyrtis puella-sinensis Ehrenberg [sic.]: Ogane et al., 2009, pl. 48, figs. 9a-9f.

Podocyrtis (Lampterium) puellasinensis Ehrenberg: Meunier and Danelian, 2023, p. 21, figs. 7.9–7.11.

Podocyrtis (Lampterium) sinuosa Ehrenberg

Podocyrtis sinuosa Ehrenberg, 1874, p. 253; 1876, p. 82, pl. 15, fig. 5.

Podocyrtis (Lampterium) sinuosa Ehrenberg: Riedel and Sanfilippo, 1970, p. 534, pl. 11, figs. 3, 4; Sanfilippo et al., 1985, p. 698, fig. 30.9.

Subgenus Podocyrtis Ehrenberg, 1847





Podocyrtis (Podocyrtis) papalis Ehrenberg, 1847

Subgenus Podocyrtoges Sanfilippo and Riedel, 1992



Podocyrtis (Podocyrtoges) ampla Ehrenberg, 1874

Podocyrtis (Podocyrtis) papalis Ehrenberg

Podocyrtis papalis Ehrenberg, 1847, p. 55, fig. 2; 1854, p. 21, pl. 36, fig. 23; 1874; p. 251; 1876, p. 82, pl. 25, fig. 6.

Podocyrtis (*Podocyrtis*) papalis Ehrenberg: Riedel and Sanfilippo, 1970, p. 533, pl. 11, fig. 1; Sanfilippo and Riedel, 1973, p. 531, pl. 20, figs. 11–14; pl. 36, figs. 2, 3.

Podocyrtis (Podocyrtoges) ampla Ehrenberg

Podocyrtis ? ampla Ehrenberg, 1874, p. 248; 1876, p. 80, pl. 16, fig. 7.

Podocyrtis (Podocyrtis) ampla Ehrenberg: Riedel and Sanfilippo,1970, p. 533, pl. 12, figs. 7, 8.

Podocyrtis (Podocyrtoges) ampla Ehrenberg: Sanfilippo and Riedel, 1992, p. 14, pl. 5, fig. 4.



Podocyrtis (Podocyrtoges) diamesa Riedel and Sanfilippo, 1970



Podocyrtis (Podocyrtoges) phyxis Sanfilippo and Riedel, 1973

Subgenus *Podocyrtopsis* Sanfilippo and Riedel, 1992



Podocyrtis (Podocyrtopsis) apeza Riedel and Sanfilippo, 1970 Not trained by the CNN

Podocyrtis (Podocyrtoges) diamesa Riedel and Sanfilippo

Podocyrtis (Podocyrtis) diamesa Riedel and Sanfilippo, 1970, p. 533 pl. 12, fig. 4 (*partim.*); Sanfilippo and Riedel, 1973, p. 531, pl. 20, figs. 9, 10, pl. 35, figs. 10, 11.

Podocyrtis (Podocyrtoges) diamesa Sanfilippo and Riedel, 1992, p. 14; Nigrini et al., 2005, p. 46, pl. P5, fig. 10.

Podocyrtis (Podocyrtoges) phyxis Sanfilippo and Riedel

Podocyrtis (Podocyrtis) diamesa Riedel and Sanfilippo, 1970, p. 533, pl. 12, fig. 6 (*partim.*). *Podocyrtis (Podocyrtis) physis* Sanfilippo and Riedel, 1973, p. 531.

Podocyrtis (Podocyrtoges) physis Sanfilippo and Riedel: Sanfilippo and Riedel, 1992, p. 14.

Podocyrtis (Podocyrtopsis) apeza Sanfilippo and Riedel

Podocyrtis (Podocyrtopsis) apeza Sanfilippo and Riedel, 1992, p. 14, pl. 3, figs. 13–15; Moore and Kamikuri, 2012, p. 10, pl. P7, fig. 7; Meunier and Danelian, 2022, p. 22, pl. 2, fig. 10.

Genus Theocyrtis Haeckel, 1887



Theocyrtis scolopax (Ehrenberg, 1874)



Theocyrtis spp.

Pterocorythiids with several tightly packed pores. Wide conical-tubular cephalis, sometimes bearing a short conical horn. Conical to campanulate thorax. Abdomen more or less complete.
Theocyrtis scolopax (Ehrenberg)

Eucyrtidium scolopax Ehrenberg, 1874, p. 232; 1876, p. 72, pl. 9, fig. 5; Ogane et al., pl. 58, figs. 3a–3f.

Theocyrtis scolopax (Ehrenberg): Popova et al., 2002, p. 50, fig. 14G.

Family **Theocotylidae** Petrushevskaya, 1981 Genus *Pterocodon* Ehrenberg, 1847



Pterocodon apis Ehrenberg, 1874 *Not trained by the CNN*

Genus Theocorys Haeckel, 1882



Theocorys spp.

Included in this class are *Theocorys anaclasta* Riedel and Sanfilippo, 1970 and *Theocorys anaphographa* Riedel and Sanfilippo, 1970

Genus *Theocotyle* Riedel and Sanfilippo, 1970



Theocotyle spp.

Included in this class are *Theocotyle venezuelensis* Riedel and Sanfilippo, 1970 and *Theocotyle cryptocephala* (Ehrenberg). 253

Pterocodon apis Ehrenberg, 1874 *Pterocodon apis* Ehrenberg, 1874, p. 255. Family **Theocotylidae** Petrushevskaya, 1981 Genus *Thyrsocyrtis* Ehrenberg, 1847 Subgenus *Pentalocorys* Haeckel, 1882





Thyrsocyrtis (Pentalocorys) krooni (Ehrenberg, 1874) *Not trained by the CNN*







Thyrsocyrtis (Pentalocorys) triacantha (Ehrenberg, 1874)

Subgenus Thyrsocyrtis Ehrenberg, 1847









Thyrsocyrtis (Thyrsocyrtis) rhizodon Ehrenberg, 1874

Thyrsocyrtis (Pentalacorys) krooni Sanfilippo and Blome

Thyrsocyrtis tetracantha (Ehrenberg): Riedel and Sanfilippo, 1978, p. 81, pl. 10, fig. 9 (partim.).

Thyrsocyrtis (Pentalacorys) tetracantha (Ehrenberg): Sanfilippo and Riedel, 1982, p. 176 (*partim.*), pl. 1, fig. 11.

Thyrsocyrtis (Pentalacorys) krooni Sanfilippo and Blome, 2001, p. 207, figs. 7a–7e; Moore and Kamikuri, 2012, p. 11, pl. P8, figs. 7, 8; Meunier and Danelian, 2022, p. 23, pl. 2, fig. 14.

Thyrsocyrtis (Pentalacorys) triacantha (Ehrenberg)

Podocyrtis triacantha Ehrenberg, 1874, p. 254; 1876, p. 82, pl. 13, fig. 4

Thyrsocyrtis (Pentalacorys) triacantha (Ehrenberg): Sanfilippo and Riedel, 1982, p. 176, pl. 1, figs. 8–10, pl. 3, figs. 3, 4.

Thyrsocyrtis (Thyrsocyrtis) rhizodon Ehrenberg

Thyrsocyrtis rhizodon Ehrenberg, 1874, p. 262; 1876, p. 84, pl. 12, fig. 1; Sanfilippo and Riedel, 1982, p. 173, pl. 1, figs. 14–16; pl. 3, figs. 12–17.

Thyrsocyrtis (Thyrsocyrtis) rhizodon Ehrenberg: Kamikuri, 2015, pl. 19, figs. 6a, 6b.

Family **Theoperidae** Haeckel, 1882 Genus *Eusyringium* Haeckel, 1882



Eusyringium fistuligerum (Ehrenberg, 1874) group

All *Eusyringium* species with a distal tube is placed here no matter if there is or not a stricture present between the thorax and the tube. Our specimens don't have any strictures but since there are a lot of different *Eusyringium* species and variations described and published by different authors we decided to refer all or *Eusyringium* with distal tubes as *E. fistuligerum* group.



Eusyringium lagena (Ehrenberg, 1874)

Eusyringium fistuligerum (Ehrenberg) group

Eucyrtidium tubulus Ehrenberg, 1854, p. 21, pl. 36, fig. 19; 1874, p. 233; 1876, p. 72, pl. 9, fig. 6.

Eusyringium tubulus (Ehrenberg): Ling, 1975, p. 729, pl. 9, fig. 22.

Eucyrtidium fistuligerum Ehrenberg, 1874, p. 229; 1876, p. 70, pl. 9, fig. 3.

Eusyringium fistuligerum (Ehrenberg): Riedel and Sanfilippo, 1970, p. 527, pl. 8, figs. 8 ?, 9; Moore, 1971, p. 741, pl. 4, fig. 10–11; Petrushevskaya and Kozlova, 1972, p. 549, pl. 32, fig. 3–5; Ling, 1975, p. 728, pl. 9, figs. 19, 20; Sanfilippo and Blome, 2001, p. 212, fig. 9a–d.

Eusyringium lagena (Ehrenberg)

Lithopera lagena Ehrenberg, 1874, p. 241; 1876, p. 78, pl. 3, fig. 4.

Eusyringium lagena (Ehrenberg): Riedel and Sanfilippo, 1970, p. 527, pl. 8, figs. 5–7; Foreman, 1973, p. 435, pl. 11, figs. 4, 5.

Genus Rhopalocanium Ehrenberg, 1847



Rhopalocanium ornatum (Ehrenberg, 1874) group

Three-segmented pterocorythoids with with two to three lateral wings. On some specimens the width of the abdomen narrows down distally to form a tube.



Rhopalocanium sphinx (Ehrenberg, 1874) *Not trained by the CNN*



Pterocorythoid group

Three-segmented Pterocorythoids, displaying an elongated campanulate cephalis and a campanulate to inflated spherical thorax. The abdomen is sometimes not fully developed, in some other specimens it is more cylindrical and tapering distally, on others it is slightly inflated. The length of the abdomen is shorter or equal to the length of the thorax. 259

Rhopalocanium ornatum Ehrenberg group

Rhopalocanium ornatum Ehrenberg, 1847, fig. 3; 1854, pl. 36, fig. 9; 1874, p. 256; 1876, p. 82, pl. 17, fig. 8; Foreman, 1973, p. 439, pl. 2, figs. 8–10, pl. 12, fig. 3; Riedel and Sanfilippo, 1978, p. 72, pl. 9, fig. 5; Sanfilippo and Blome, 2001, p. 217, figs. 10o, 10p.

Remarks: There are a lot of morphological variabilities within this group which leads us to determine this as a group of species.

Rhopalocanium sphinx (Ehrenberg, 1874)

Rhopalocanium sphinx (Ehrenberg, 1874); p. 255.



Nassellaria group A A trash group of nassellarians.



Nassellaria group B

Pterocorythids with a large cephalis. Pores of the thorax are similar to those of the abdomen. Many specimens have a long abdomen, no horn and a discrete outline (weak collar stricture) between the thorax and cephalis.



Nassellaria group C

50 µm

Small nassellarians with two to three segments that bear two lateral wings expanding from the thorax. The abdomen of some specimens is very long and cylindrical. *Not trained by the CNN.*

Spumellarians

Superfamily LITHOCYCLIOIDEA Ehrenberg, 1846 Family Lithocycliidae Ehrenberg, 1846 Genus Lithocyclia Ehrenberg, 1847



Lithocyclia ocellus Ehrenberg, 1854 group

Family **Phacodiscidae** Haeckel, 1882 Genus *Periphaena* Ehrenberg, 1874





Periphaena decora Ehrenberg, 1874 Some of the encountered forms display small equatorial spines..

Lithocyclia ocellus group

Lithocyclia ocellus Ehrenberg group

Lithocyclia ocellus Ehrenberg, 1854b, pl. 36, fig. 30; Petrushevskaya and Kozlova, 1972, p. 523, pl. 15, figs. 1-2.

Lithocyclia ocellus Ehrenberg group: Riedel and Sanfilippo, 1970, p. 522, pl. 5, figs. 1-2; Riedel and Sanfilippo, 1971, p. 1588, pl. 3A, fig. 6; Kamikuri, 2015, pl. 19, fig. 2.

Periphaena decora Ehrenberg

Periphaena decora Ehrenberg, 1874, p. 246; 1876, p. 80, pl. 28, fig. 6; Petrushevskaya and Kozlova, 1972, p. 523, pl. 14, figs. 1, 2; Sanfilippo and Riedel, 1973, p. 523, pl. 8, figs. 8-10, pl. 27, figs. 2-4 (partim.); Nigrini, 1974, p. 1065, pl. 1C, figs. 1, 2, 4, 6 (partim.); Funakawa et al., 2006, p. 42, pl. P15, figs. 5a-6b; Ogane et al., 2009, pl. 14, figs. 6a-6c, pl. 30, fig. 4, pl. 67, figs. 3a, 3b.

Superfamily LITHELIOIDEA Haeckel, 1862 Family Litheliidae Haeckel, 1862 Genus *Middourium* Kozlova, 1999





"Middourium" group

We include in this class forms that are characterized by an ellipsoidal, loosely concentric shell that resembles *Middourium regulare* sensu O'Dogherty et al. (2021, p. 873). However, some of these forms may be incomplete or broken shells of *Lithelius* spp.

Superfamily **STYLOSPHAEROIDEA** Haeckel, 1887 Family **Stylosphaeroidae** Haeckel, 1887

Genus Spongatractus Hackel, 1887



Spongatractus pachystylus (Ehrenberg, 1874) Not trained by the CNN

Genus Stylosphaera Ehrenberg, 1846



Stylosphaera coronata Ehrenberg, 1874 group

Single shelled spherical stylosphaeroids, most of which bear two polar spines of variable length and thickness. The typical *Stylosphaera coronata* have two spines, while we also encountered forms which have four spines which is why we consider this to be a possibility of being a group of species. 267

Stylosphaera coronata Ehrenberg, 1874

Stylosphaera coronata Ehrenberg, 1874, p. 258.

Remarks: There are a lot of morphological variabilities within this group which leads us to determine this as a group of species.

Spongatractus pachystylus (Ehrenberg)

Spongatractus pachystylus (Ehrenberg): Sanfilippo and Riedel, 1973, p. 519, pl. 2, figs. 4-6, pl. 25, fig. 3.

Spongosphaera pachystyla Ehrenberg, 1874, p. 256; 1876, p. 82, pl. 26, fig. 3.

Superfamily HELIOSATURNALOIDEA Kozur & Mostler, 1972 Family Axoprunidae Dumitrica, 1985

Genus Axoprunum Haeckel 1887



"Axoprunum" sp. A

Simple spherical spumellarians with large pores having two polar long elongated spiny spines.

Not trained by the CNN





Discoidal biconvex spumellarians with equatorial spines

In this class are included forms that belong to at least two genera: *Heliodiscus* (family Heliodiscidae) and *Periphaena* (family Phocodicidae).



Spherical spumellarians with large pores

Sphaerical spumellarians (some of them multishelled) displaying large subcircular pores on the outer shell.



Multi-armed spongy spumellarians

In this class are included spongy spumellarians composed of a variable number of arms and a patagium developed in their central part. Spines may be present or absent. One genus included in this class is *Histiastrum* Ehrenberg 1847.





Spherical spumellarians with radial spines

Sphaerical multishelled or spongy spumellarians with numerous (more than 8) conical rodlike spines.

Not trained by the CNN



Discoidal spongy spumellarians with rod shaped spinesWe included in this class flat discoidal spongy shells bearing a few elongated rod-shaped spines.



Discoidal spongy spumellarians

We have included in this class discoidal forms that are composed internally of a shell that appears usually spongy or sometimes displays concentric rings. 271



Discoidal spumellarians with regular concentric rings

In this class are included forms that belong to at least two families Trematodiscidae, such as *Frustrella* and Circodiscidae. They either display short spines or no spines.



Discoidal spumellarians with *irregular* concentric rings

In this class are included forms that belong to at least two families Trematodiscidae displaying irregular concetric rings.





Small discoidal spumellarians with concentric rings

In this class are included forms that belong to at least the family Trematodiscidae displaying concetric rings. Some specimens can be broken central parts of larger Trematodisciids.

References

Bjørklund, K. R. (1976). Radiolaria from the Norwegian Sea, Leg 38 of the Deep Sea Drilling Project. Initial Reports of the Deep Sea Drilling Project, 38, 1101-1168.

Blueford, J. (1988). Radiolarian biostratigraphy of siliceous Eocene deposits in central California. Micropaleontology, 34(3), 236-258. doi:10.2307/1485754.

Bütschli, O. (1882). Beiträge zur Kenntnis des Radiolarienskelettes, insbesondere der der Cyrtida. Zeitschrift für wissenschaftliche Zoologie, 36, 486-540.

Campbell, A. S., & Clark, B. L. (1944). Radiolaria from Upper Cretaceous of middle California. Geological Society of America, Special Papers, 57, 1-61.

Clark, B. L., & Campbell, A. S. (1942). Eocene radiolarian faunas from the Mt. Diablo area, California. Geological Society of America, Special Papers, 39, 1-112.

De Wever, P., Dumitrica, P., Caulet, J. P., Nigrini, C., & Caridroit, M. (2001). Radiolarians in the sedimentary record. London, UK: Gordon & Breach Science Publishers, 533.

Dumitrica, P. (1985). Internal morphology of the Saturnalidae (Radiolaria): systematic and phylogenetic consequences. Revue de Micropaleontologie, 28(3), 181-196.

Ehrenberg, C. G. (1846). Über eine halibiolithische, von Herrn R. Schomburgk entdeckte, vorherrschend aus mikroskopischen Polycystinen gebildete, Gebirgsmasse von Barbados. Bericht über die zur Bekanntmachung geeigneten Verhandlungen der Königlich Preussischen Akademie der Wissenschaften zu Berlin, 382-385.

Ehrenberg, C. G. (1847). Über die mikroskopischen kieselschaligen Polycystinen als machtige Gebirgsmasse von Barbados und über das Verhältniss der aus mehr als 300 neuen Arten bestehenden ganz eigenthümlichen Formengruppe jener Felsmasse zu den jetzt lebenden Tieren und zur. Bericht über die zu Bekanntmachung geeigneten Verhandlungen der Königlichen Preussische Akademie der Wissenschaften zu Berlin, 41-60.

Ehrenberg, C. G. (1854). Mikrogeologie. Leipzig, 1-374.

Ehrenberg, C. G. (1875). Fortsetzung der mikrogeologischen Studien. Abhandlungen der königlichen Akademie der Wissenschaften (Berlin).

Ehrenberg, C. G. (1876). Fortsetzung der mikrogeologischen Studien als Gesammt-Übersicht der mikroskopischen Paläontologie gleichartig analysirter Gebirgsarten der Erde, mit specieller Rücksicht auf den Polycystinen-Mergels von Barbados. Königlich Preussischen Akademie der Wissenschaften zu Berlin, Abhandlungen, Jahre 1875, 1-226.Ehrenberg, C. G. 1874. Grössere Felsproben des Polycystinen-Mergels von Barbados mit weiteren Erläuterungen. Monatsberichte der Königlichen Preussischen Akademie der Wissenschaften zu Berlin 213-263.

Ehrenberg, C. G. (1875). Fortsetzung der mikrogeologischen Studien. Abhandlungen der königlichen Akademie der Wissenschaften (Berlin).

Ehrenberg, C. G. (1876). Fortsetzung der mikrogeologischen Studien als Gesammt-Übersicht der mikroskopischen Paläontologie gleichartig analysirter Gebirgsarten der Erde, mit specieller Rücksicht auf den Polycystinen-Mergels von Barbados. Königlich Preussischen Akademie der Wissenschaften zu Berlin, Abhandlungen, Jahre 1875, 1-226.

Empson-Morin, K. M. (1981). Campanian Radiolaria from DSDP Site 313, Mid-Pacific Mountains. Micropaleontology, 27(3), 249-292. 273 Foreman, H. P. (1973). Radiolaria of Leg 10 with systematics and ranges for the families Amphipyndacidae, Artostrobiidae, and Theoperidae. Initial Reports of the Deep Sea Drilling Project, 10, 407-474.

Funakawa, S. (1994). Plagiacanthidae (Radiolaria) from the Upper Miocene of eastern Hokkaido, Japan. Transactions and Proceedings of the Palaeontological Society of Japan, 174, 458-483.

Funakawa, S., Nishi, H., Moore, T. C., Jr., and Nigrini, C. A. (2006). Data report: Late Eocene early Oligocene radiolarians, ODP Leg 199 Holes 1218A, 1219A, and 1220A, central Pacific. In Wilson, P. A., Lyle, M., Firth, J. V. (Eds.), Proceedings of the Ocean Drilling Program, Scientific Results, 199. Ocean Drilling Program, College Station, TX, 1-74. doi:10.2973/odp.proc.sr.199.216.2006.Goll, R. M. (1968). Classification and phylogeny of Cenozoic Trissocyclidae (Radiolaria) in the Pacific and Caribbean basins, Part I. Journal of Paleontology, 42(6), 1409-1432.

Goll, R. M. (1969). Classification and phylogeny of Cenozoic Trissocyclidae (Radiolaria) in the Pacific and Caribbean basins, Part II. Journal of Paleontology, 43(2), 322-339.

Haeckel, E. 1862. Die Radiolarien (Rhizopoda Radiaria). Eine Monographie. Reimer, Berlin, 572 p. doi:10.5962/bhl.title.10155

Haeckel, E. 1882. Entwurf eines Radiolarien-Systems auf Grund von Studien der Challenger-Radiolarien. Jenaische Zeitschrift für Naturwissenschaft herausgegeben von der medizinisch-naturwissenschaftlichen Gesellschaft zu Jena 15: 418-472.

Haeckel, Ernst. 1887. Report on the Radiolaria collected by H.M.S. Challenger during the years 1873-1876. Report on the Scientific Results of the Voyage of H.M.S. Challenger during the years 1873-1876. 18: 1-1803.

Haecker, V. (1908). Tiefsee-Radiolarien. Spezieller Teil. Die Tripyleen, Collodarien und Mikroradiolarien der Tiefsee. In: Chun, C. (Ed.), Wissenschaftliche Ergebnisse der Deutschen Tiefsee-Expedition auf dem Dampfer Valdivia, 1898-1899, Vol. 14, Jena, Germany, 336-476.

Hertwig, R. (1879). Der Organismus der Radiolarien. G. Fischer, Jena, Germany, iv + 149 p. https://archive.org/details/denkschriftender02medi/page/126/mode/2up

Hollis, C. J., Pascher, K. M., Sanfilippo, A., Nishimura, A., Kamikuri, S.-i., and Shepherd, C. L. (2020). An Austral radiolarian biozonation for the Paleogene. Stratigraphy, 17(4), 213-278. doi: 10.29041/strat.17.4.213-278.

Hull, D. M. (1996). Paleoceanography and biostratigraphy of Paleogene radiolarians from the Norwegian-Greenland Sea. In: Thiede, J., Myhre, A. M., Firth, J. V., Johnson, G. L., and Riddiman, W. F. (Eds.), Proceedings of the Ocean Drilling Program, Scientific Results, 151, 125-152. Ocean Drill. Program, College Station, TX, USA. doi:10.2973/odp.proc.sr.151.103.1996.

Johnson, D.A. (1974). Radiolaria from the Eastern Indian Ocean, DSDP Leg 22, in: Von der Borch, C.C., Sclater, J.G., Gartner, Jr., S., Hekinian, R., Johnson, D.A., McGowran, B., Pimm, A.C., Thompson, R.W., Veevers, J.J., Waterman, L.S. (eds), Initial Reports of the Deep Sea Drilling Project, Volume 22. U.S. Government Printing Office, Washington, 521-575. doi:10.2973/DSDP.PROC.22.125.1974.

Kamikuri, S. (2015). Radiolarian assemblages during the middle-late Eocene transition at Site 1052, ODP Leg 171B, Blake Nose, western North Atlantic Ocean. News of Osaka Micropaleontologists, Special Volume 15, 139-167.

Kim, K. Hyune. 1992. Paleogene radiolarian biostratigraphy from high-latitude South Atlantic. Journal of the Paleontological Society of Korea 8(1): 24-51.

Kozur, H., & Mostler, H. (1972). Beiträge zur Erforschung der mesozoischen Radiolarien. Teil I: Revision der Oberfamilie Coccodiscacea Haeckel 1862 emend. und Beschreibung ihrer triassischen Vertreter. Geologisch Palaontologische Mitteilungen Innsbruck, 2, 1-60.

Kozlova, G. E. (1999). Radiolyarii Paleogena Boreal'noy Oblasti Rossii. [Paleogene Boreal Radiolarians from Russia.] Prakticheskoe Rukovodstvo po Mikrofaune Possii 9. 1-323.

Lazarus, D. B., & Pallant, A. (1989). Oligocene and Neogene radiolarians from the Labrador Sea: ODP Leg 105. In: Srivastava, S. P., Arthur, M., Clement, B., et al. (Eds.), Proceedings of the Ocean Drilling Program, Scientific Results, 105, 349-380. College Station, TX: Ocean Drilling Program.

Ling, H.Y. (1975). Radiolaria: Leg 31 of the Deep Sea Drilling Project, in: Karig, D.E., Ingle, J.C., Jr., Bouma, A.H., Ellis, C.H., Haile, N., Koizumi, I., Ling, H.Y., MacGregor, I., Moore, J.C., Ujiie, H., Watanabe, T., White, S.M., Yasui, M. (eds), Initial Reports of the Deep Sea Drilling Project 31. U.S. Government Printing Office, Washington, 703-761.

Meunier, M., and Danelian, T. (2022). Astronomical calibration of late middle Eocene radiolarian bioevents from ODP Site 1260 (equatorial Atlantic, Leg 207) and refinement of the global tropical radiolarian biozonation. Journal of Micropalaeontology, 41(1), 1-27.

Meunier, M., & Danelian, T. (2023). Progress in understanding middle Eocene nassellarian (Radiolaria, Polycystinea) diversity; new insights from the western equatorial Atlantic Ocean. Journal of Paleontology, 97(1), 1-25. doi:10.1017/jpa.2022.82

Moore Jr., T.C. (1971). Radiolaria. In: Initial Reports of the Deep Sea Drilling Project, Volume 8, edited by J.I. Tracey Jr., G.H. Sutton, W.D. Nesteroff, J. Galehouse, C.C. Von der Borch, T. Moore, J. Lipps, U.Z.B.U. Haq, and J.P. Beckmann, U.S. Govt. Print. Office, Washington, DC, USA, pp. 727-775. doi:10.2973/dsdp.proc.8.112.1971.

Nigrini, C. (1974). Cenozoic Radiolaria from the Arabian Sea, DSDP Leg 23. In: Whitmarsh, R.B., Weser, O.E., Ali, S., Boudreaux, J.E., Fleisher, R.L., Jipa, D., Kidd, R.B., Mallik, T.K., Matter, A., Nigrini, C., Siddiquie, H.N., Stoffers, P. (eds.), Initial Reports of the Deep Sea Drilling Project, Volume 23, U.S. Govt. Print. Office, Washington, pp. 1051-1121. doi:10.2973/DSDP.PROC.26.233.1974.

Nigrini, C., (1977). Equatorial Cenozoic Artostrobiidae (Radiolaria). Micropaleontology 23 (3), 241–269.

Nigrini, C., Sanfilippo, A., Moore, T.J. Jr. (2005). Cenozoic radiolarian biostratigraphy: a magnetobiostratigraphic chronology of Cenozoic sequences from ODP Sites 1218, 1219, and 1220, equatorial Pacific. In Wilson, P.A., Lyle, M., Firth, J.V. (Eds.), Proc. ODP, Sci. Results 199. Ocean Drill. Program, College Station, TX, pp. 1–76. doi:10.2973/odp.proc.sr.199.225.2005.

Nigrini, C. A., Sanfilippo, A. & Moore, T. C. Jr. (2006). Cenozoic radiolarian biostratigraphy: a magnetobiostratigraphic chronology of Cenozoic sequences from ODP Sites 1218, 1219, and 1220, equatorial Pacific. Proceedings of the Ocean Drilling Program, Scientific Results. 199: 1-76.

O'Connor, B. (1999b). Radiolaria from the Late Eocene Oamaru Diatomite, South Island, New Zealand. Micropaleontology, 45(1), 1-55. doi:10.2307/1486169

Ogane, K., Suzuki, N., Aita, Y., Sakai, T., Lazarus, D., Tanimura, Y. (2009). Ehrenberg's radiolarian collections from Barbados. In Tanimura, Y., Aita, Y. (Eds.), Joint Haeckel and Ehrenberg Project: Reexamination of the Haeckel and Ehrenberg Microfossil Collection as a historical and scientific legacy. National Museum of Nature and Science Monographs, Tokyo, 40, 97-106.

Petrushevskaya, M. G. (1967). Radiolayarii otryadov Spumellaria i Nassellaria Antarkicheskoi oblasti (po materialam Sovetskoi Antarkitcheskikh Ekspeditsii) [Radiolarians of the orders Spumellaria and Nassellaria from the Antarctic region (based on the material of the Soviet Antarctic Expedition)]. Issledovanie Faunyi Morey (Akademii Nauk SSSR) Leningrad, 4(12), 5-186. [in Russian]

Petrushevskaya, M. G., and Kozlova, G. E. (1972). Radiolaria: Leg 14, Deep Sea Drilling Project. Initial Reports of the Deep Sea Drilling Project, 14, 459-648.

Petrushevskaya, M.G., Kozlova, G.E. (1979). Opisanie rodov i vidov Radiolyarii. Issledovaniya Fauny Morei, 23, 86-157 [in Russian].

Petrushevskaya, M. G. (1981). Nassellarian radiolarians from the world oceans. Publications of the Zoological Institute, Academy of Sciences of the USSR, Nauka, Leningradskoe Otdelenie, Leningrad, USSR, 405 p. [in Russian]

Pessagno Jr, E. A. (1976). Radiolarian zonation and stratigraphy of the Upper Cretaceous portion of the Great Valley sequence, California Coast Ranges. Micropaleontology, Special Publication 2, 1-95.

Popova, I.M., Baumgartner, P.O., Guex, J., Tochilina, S.V., Glezer, Z.I. (2002). Radiolarian biostratigraphy of Palaeogene deposits of the Russian Platform (Voronesh Anticline). Geodiversitas, 24(1), 7-59.

Renz, G. W. (1984). Cenozoic radiolarians from the Barbados Ridge, Lesser Antilles subduction complex, Deep Sea Drilling Project Leg 78A. Initial Reports of the Deep Sea Drilling Project, 78A, 447-462.

Riedel, W.R. (1957). Radiolaria: a preliminary stratigraphy. In Petterson, H. (Ed.), Reports of the Swedish Deep-Sea Expedition, 1947-1948. Elanders Boktryckeri Aktiebolag, Göteborg, Sweden, 6, 59-96.

Riedel, W. R. (1967b). Some new families of Radiolaria. Proceedings of the Geological Society of London, 1640, 148-149.

Riedel, W. R., & Sanfilippo, A. (1970). Radiolaria, Leg. 4 DSDP. Initial Reports of the Deep Sea Drilling Project, 4, 503-575.

Riedel, W. R., & Sanfilippo, A. (1971). Cenozoic Radiolaria from the western tropical Pacific, Leg 7. Initial Reports of the Deep Sea Drilling Project, 7, 1529-1666.

Riedel, W.R., Sanfilippo, A. (1977). Cenozoic Radiolaria. In Ramsay, A.T.S. (Ed.), Oceanic Micropalaeontology, Vol. 2. Academic Press, London/New York/San Francisco, 847-912.

Riedel, W.R., Sanfilippo, A. (1978). Stratigraphy and evolution of tropical Cenozoic radiolarians. Micropaleontology, 24, 61-96. doi:10.2307/1485420.

Salamé, L., & Huber, B. T. (2014). Biostratigraphy of Late Paleocene - Middle Eocene radiolarians and foraminifera from Cyprus. Micropaleontology, 60(3-4), 237-261. doi:10.2113/gsmicropal.60.3.237

Sandin, M. M., Biard, T., Romac, S., O'Dogherty, L., Suzuki, N., Not, F. (2021). A Morpho-molecular Perspective on the Diversity and Evolution of Spumellaria (Radiolaria). Protist, 172(3), 125806. ISSN 1434-4610. doi:10.1016/j.protis.2021.125806.

Sanfilippo, A., & Riedel, W. R. (1973). Cenozoic Radiolaria (exclusive of theoperids, artostrobiids and amphipyndacids) from the Gulf of Mexico, Deep Sea Drilling Project Leg 10. Initial Reports of the Deep Sea Drilling Project, 10, 475-608.

Sanfilippo, A., Caulet, J.-P. (1998). Taxonomy and evolution of Paleogene Antarctic and tropical Lophocyrtid radiolarians. Micropaleontology, 44(1), 1-43.

Sanfilippo, A. & Riedel, W. R. (1979). Radiolaria from the northeastern Atlantic Ocean DSDP Leg 48. Initial Reports of the Deep Sea Drilling Project. 48: 493-511.

Sanfilippo, A., Riedel, W.R. (1982). Revision of the radiolarian genera Theocotyle, Theocotylissa and Thyrsocyrtis. Micropaleontology, 28(2), 170-188.

Sanfilippo, A., Westberg-Smith, M.J., Riedel, W.R. (1985). Cenozoic Radiolaria. In Bolli, H.M., Saunders, J.B., Perch-Nielsen, K. (Eds.), Plankton Stratigraphy. Cambridge University Press, Cambridge, UK, 631-712.

Sanfilippo, A. (1990). Origin of the subgenera Cyclampterium, Parampterium and Sciadiopeplus from Lophocyrtis (Lophocyrtis) (Radiolaria, Theoperidae). Marine Micropaleontology, 15(3-4), 287-312. doi:10.1016/0377-8398(90)90024-A

Sanfilippo, A., & Riedel, W. R. (1992). The origin and evolution of Pterocorythidae (Radiolaria): A Cenozoic phylogenetic study. Micropaleontology, 38, 1-36.

Sanfilippo, A., & Caulet, J. P. (1998). Taxonomy and evolution of Paleogene Antarctic and Tropical Lophocyrtid radiolarians. Micropaleontology, 44(1).

Sanfilippo, A., Blome, C.D. (2001). Biostratigraphic implications of mid-latitude Palaeocene-Eocene radiolarian faunas from Hole 1051A, ODP Leg 171B, Blake Nose, western North Atlantic. In Kroon, D., Norris, R.D., Klaus, A. (Eds.), Western North Atlantic Palaeogene and Cretaceous Palaeoceanography. Geological Society Special Publication, 183(1), 185-224. doi:10.1144/GSL.SP.2001.183.01.10.

Schröder, O. (1909). Die nordischen Spumellarian: Unterlegion Sphaerellaria. In Brandt, K. & Apstein, C. (Eds.), Nordisches Plankton (Vol. 17, pp. 1-66). Lipsius und Tischer.

Shilov, V.V. (1995). Eocene-Oligocene Radiolarians from Leg 145, North Pacific. In Rea, D.K., Basov, I.A., Scholl, D.W., Allan, J.F. (Eds.), Proceedings of the Ocean Drilling Program, Scientific Results, 145, 117-132.

Suzuki, N., Ogane, K., Chiba, K. (2009). Middle to Late Eocene polycystine radiolarians from Site 1172, Leg 189, Southwest Pacific. News of Osaka Micropaleontologists, Special Volume 14, 239-296.

Suzuki, Noritoshi, Caulet, Jean-Pierre, & Dumitrica, Paulian. (2021). Cycladophoridae Suzuki 2019. In A new integrated morpho- and molecular systematic classification of Cenozoic radiolarians (Class Polycystinea) - suprageneric taxonomy and logical nomenclatorial acts, pp. 405-573 in Geodiversitas (Vol. 43, Number 15, pp. 513–514). Zenodo. doi:10.5281/zenodo.5106785

Suzuki, N., O'Dogherty, L., Caulet, J.-P., & Dumitrica, P. (2021). A new integrated morpho- and molecular systematic classification of Cenozoic radiolarians (Class Polycystinea) – suprageneric taxonomy and logical nomenclatorial acts. Geodiversitas, 43(15), 405-573.

Takemura, A. (1992). Radiolarian Paleogene biostratigraphy in the southern Indian Ocean, Leg 120. In Schlich, R., Wise, S.W., Jr., Palmer Julson, A.A., Aubry, M.-P., Berggren, W.A., Bitschene, P.R., Blackburn, N.A., Breza, J., Coffin, M.F., Harwood, D.M., Heider, F., Holmes, M.A., Howard, W.R., Inokuchi, H., Kelts, K., Lazarus, D.B., Mackensen, A., Muruyama, T., Munschy, M., Pratson, E., Quilty, P.G., Rack, F., Salters, V.J.M., Sevigny, J.H., Storey, M., Takemura, A., Watkins, D.K., Whitechurch, H., Zachos, J. (Eds.), Proc. ODP, Sci. Results 120. Ocean Drill. Program, College Station, TX, pp. 735-756. doi:10.2973/odp.proc.sr.120.177.1992.

Zittel, K. A. (1876). Über einige fossile Radiolarien aus der norddeutschen Kreide. Zeitschrift der Deutschen Geologischen Gesellschaft, 28, 75-86.

Abstract

Micropaleontology is not only about studying organisms that lived in the past, with implications for biostratigraphy and evolutionary changes within morphospecies. it is also about understanding Earth's past environments, paleoceanography and climate change. This field is facing numerous challenges, since the analysis of microfossils demands significant human effort and taxonomic expertise. This PhD work focuses on the application of Artificial Intelligence (AI), such as Artificial Neural Networks (ANNs), for automatic image recognition of tropical Atlantic middle Eocene radiolarians. Large datasets have been constructed, in order to train different neural networks. Our results show that neural networks can automatically classify several different radiolarian classes down to the species level; in many cases, they were able to identify closely related species and even transitional morphotypes in evolutionary lineages. It was also possible to correctly identify partly broken or blurry radiolarians. NNs were also successfully applied to automatic image recognition to be used with a biostratigraphic aim. This work includes the use of the classical neural network approaches for analysing visual context such as Convolutional Neural Networks (CNNs), but also the use of Spiking Neural Networks (SNNs), which are not as commonly used for automatic image recognition, as CNNs. SNNs resulted in nearly equal accuracy as the one obtained with CNNs, but the use of an SNN is more computational efficient and takes up less memory. There has also been a comparison between the outcome of NN training and of traditional morphometric analyses, such as Linear Discrimination Analysis (LDA), giving approximately similar results. Our research does not only simplify and speed up the analysis process, but also helps in increasing the accuracy and consistency of micropaleontological interpretations.

Keywords: Artificial Intelligence, Radiolaria, Neural Networks, middle Eocene, CNN, SNN, biostratigraphy

Résumé

La micropaléontologie n'est pas seulement concernée par l'étude d'organismes qui ont vécu dans le passé, avec des implications pour la biostratigraphie et les changements évolutifs au sein des morphoespèces, mais également par la compréhension des environnements passés de la Terre, la paléoocéanographie et les changements climatiques. Ce domaine est confronté à de nombreux défis, car l'analyse des microfossiles nécessite un effort humain et une expertise taxonomique considérables. La présente thèse traite de l'application de l'intelligence artificielle (AI), tels les réseaux de neurones artificiels (RNA), pour la reconnaissance automatique d'images de radiolaires de l'Éocène moyen de l'Atlantique tropical. Des larges bases de données ont été construites afin d'entrainer différents réseaux de neurones. Nos résultats montrent que ceux-ci réussissent à classer automatiquement plusieurs classes différentes de radiolaires, jusqu'au niveau de l'espèce, ainsi que dans de nombreux cas, d'identifier des espèces étroitement apparentées et même des morphotypes transitionnels au sein des lignées évolutives. Le RNA a même pu déterminer correctement des radiolaires brisés ou des images floues. Un RNA a également pu être appliqué avec succès à la reconnaissance automatique d'images pour un travail biostratigraphique. Ce travail inclut l'utilisation des approches classiques de réseaux de neurones pour analyser le contexte visuel, tels que les réseaux neuronaux convolutifs (CNN), mais comprend également l'utilisation de réseaux de neurones à pointes (SNN), qui ne sont pas aussi couramment utilisés que les CNN pour la reconnaissance automatique d'images. Les SNN ont permis d'obtenir une précision quasi-équivalente à celle des CNN, avec la différence que leur utilisation est plus efficace en termes de calcul et a besoin de moins de mémoire. Il y a également eu comparaisons des résultats issus des RNA et d'analyses morphométrique classiques, telles que l'analyse de discrimination linéaire (LDA), avec des résultats sensiblement similaires. Notre recherche n'a pas seulement simplifié et accéléré le processus analytique, mais a également contribué à accroître la précision et la cohérence des interprétations micropaléontologiques.

Mots clés: Intelligence artificielle, radiolaires, réseaux de neurones, Éocène moyen, CNN, SNN, biostratigraphie