# Activity Report 2018

# Team PLEIADE

# Patterns of diversity and networks of functions

Inria teams are typically groups of researchers working on the definition of a common project, and objectives, with the goal to arrive at the creation of a project-team. Such project-teams may include other partners (universities or research institutions).

# Table of contents

<div align="center">**Team PLEIADE**</div>

*Creation of the Team: 2015 January 01, updated into Project-Team: 2019 March 01*

**Keywords:**

### Computer Science and Digital Science:

        A3.1. - Data
        A3.2. - Knowledge
        A3.3.2. - Data mining
        A3.3.3. - Big data analysis
        A3.4. - Machine learning and statistics
        A6.2.8. - Computational geometry and meshes

### Other Research Topics and Application Domains:

        B1.1.7. - Bioinformatics
        B1.1.10. - Systems and synthetic biology
        B3. - Environment and planet

# 1. Team, Visitors, External Collaborators

**Research Scientists**
    David Sherman [Team leader, Inria, Senior Researcher, HDR]
    Pascal Durrens [CNRS, Researcher, HDR]
    Alain Franc [INRA, Senior Researcher, HDR]

**Post-Doctoral Fellow**
    Swarna Kanchan [Inria, from Sep 2018]

**Technical staff**
    Philippe Chaumeil [INRA]
    Jean-Marc Frigerio [INRA]
    Franck Salin [INRA]
    Louise Amelie Schmitt [Inria, until Mar 2018]

**Administrative Assistants**
    Catherine Cattaert Megrat [Inria, from Sep 2018]
    Nathalie Robin [Inria, until Aug 2018]

# 2. Overall Objectives

## 2.1. Overall Objectives

Diversity, evolution, and inheritance form the heart of modern biological thought. Modeling the complexity of biological systems has been a challenge of theoretical biology for over a century [29] and flourished with the evolution of data for describing biological diversity, most recently with the transformative development of high-throughput sequencing. However, most concepts and tools in ecology and population genetics for capitalizing on this wealth of data are still not adapted to high throughput data production. A better connection between high-throughput data production and tool evolution is highly needed: *computational biodiversity*.

Paradoxically, diversity emphasizes differences between biological objects, while modeling aims at unifying them under a common framework. This means that there is a limit beyond which some components of diversity cannot be mastered by modeling. We need efficient methods for recognizing patterns in diversity, and linking them to patterns in function. It is important to realize that diversity in function is not the same as coupling observed diversity with function. Diversity informs both the study of traits, and the study of biological functions (Figure 1). The double challenge is to measure these links quickly and precisely with pattern recognition, and to explore the relations between diversity in traits and diversity in function through modeling



*Figure 1. Diversity informs both the study of traits, and the study of biological functions*

PLEIADE links recognition of patterns, classes, and interactions with applications in biodiversity studies and biotechnology. We develop distance methods for NGS datasets at different levels of organization: between genomes, between individual organisms, and between communities; and develop high-performance pattern recognition and statistical learning techniques for analyzing the resulting point clouds. We refine inferential methods for building hierarchical models of networks of cellular functions, exploiting the mathematical relations that are revealed by large-scale comparison of related genomes and their models. We combine these methods into integrated e-Science solutions to place these tools directly in the hands of biologists.

Our methodology (Figure 2) is designed pragmatically to advance the state of the art in applications from biodiversity and biotechnology: molecular based systematics and community ecology, annotation and modeling for biotechnology.



*Figure 2.* PLEIADE *is a pluridisciplinary team. Each application in biodiversity and biotechnology follows a path calling on methods from biology (blue), mathematics (green), and computer science (red).*

# 3. Research Program

## 3.1. A Geometric View of Diversity

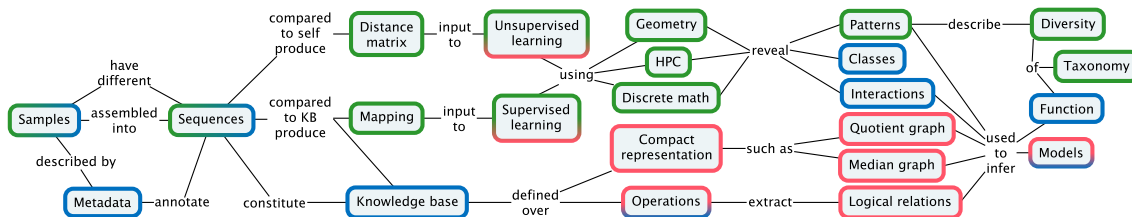Diversity may be studied as a set of dissimilarities between objects. The underlying mathematical construction is the notion of distance. Knowing a set of objects, it is possible, after computation of pairwise distances, or sometimes dissimilarities, to build a Euclidean image of it as a point cloud in a space of relevant dimension. Then, diversity can be associated with the shape of the point cloud. The human eye is often far better than an algorithm at recognizing a pattern or shape. One objective of our project is to narrow the gap between the story that a human eye can tell, and that an algorithm can tell. Several directions will be explored. First, this requires mastering classical tools in dimension reduction, mainly algebraic tools (PCA, NGS, Isomap, eigenmaps, etc ...). Second, neighborhoods in point clouds naturally lead to graphs describing the neighborhood networks. There is a natural link between modular structures in distance arrays and communities on graphs. Third, points (representing, say, DNA sequences) are samples of diversity. Dimension reduction may show that they live on a given manifold. This leads to geometry (differential or Riemannian geometry). It is expected that some properties of the manifold can tell something of the constraints on the space where measured individuals live. The connection between Riemannian geometry and graphs, where weighted graphs are seen as mesh embedded in a manifold, is currently an active field of research [28], [27]. See as well [30] for a link between geometric structure, linear and nonlinear dimensionality reduction.

Biodiversity and high-performance computing: Most methods and tools for characterizing diversity have been designed for datasets that can be analyzed on a laptop, but NGS datasets produced for metabarcoding are far too large. Data analysis algorithms and tools must be revisited and scaled up. We will mobilize both distributed algorithms like the Arnoldi method and new algorithms, like random projection or column selection methods, to build point clouds in Euclidean spaces from massive data sets, and thus to overcome the cubic complexity of computation of eigenvectors and eigenvalues of very large dense matrices. We will also link distance geometry [22] with convex optimization procedures through matrix completion [16], [17].

Intercalibration: There is a considerable difference between supervised and unsupervised clustering: in supervised clustering, the result for an item $i$ is independent from the result for an item $j \neq i$, whereas in unsupervised clustering, the result for an item $i$ (e.g. the cluster it belongs to, and its composition) depends on nearby items $j \neq i$. Which means that the result may change if some items are added to or subtracted from the sample. This raises the more global problem of how to merge two studies to yield a more comprehensive view of biodiversity?

## 3.2. Knowledge Management for Biology

The heterogenous data generated in computational molecular biology and ecology are distinguished not only by their volume, but by the richness of the many levels of interpretation that biologists create. The same nucleic acid sequence can be seen as a molecule with a structure, a sequence of base pairs, a collection of genes, an allele, or a molecular fingerprint. To extract the maximum benefit from this treasure trove we must organize the knowledge in ways that facilitate extraction, analysis, and inference. Our focus has been on the efficient representation of relations between biological objects and operations on those representations, in particular heuristic analyses and logical inference.

PLEIADE will develop applications in comparative genomics of related organisms, using new mathematical tools for representing compactly, at different scales of difference, comparisons between related genomes. New methods based on distance geometry will refine these comparisons. Compact representations can be stored, exchanged, and combined. They will form the basis of new simultaneous genome annotation methods, linked directly to abductive inference methods for building functional models of the organisms and their communities.

Since a goal of PLEIADE is to integrate diversity throughout the analysis process, it is necessary to incorporate **diversity as a form of knowledge** that can be stored in a knowledge base. Diversity can be represented using various compact representations, such as trees and quotient graphs storing nested sets of relations. Extracting structured representations and logical relations from integrated knowledge bases (Figure 2) will require domain-specific query methods that can express forms of diversity.

## 3.3. Modeling by successive refinement

Describing the links between diversity in traits and diversity in function will require comprehensive models, assembled from and refining existing models. A recurring difficulty in building comprehensive models of biological systems is that accurate models for subsystems are built using different formalisms and simulation techniques, and hand-tuned models tend to be so focused in scope that it is difficult to repurpose them [14]. Our belief is that a sustainable effort in building efficient behavioral models must proceed incrementally, rather than by modeling individual processes *de novo*. *Hierarchical modeling* [11] is one way of combining specific models into networks. Effective use of hierarchical models requires both formal definition of the semantics of such composition, and efficient simulation tools for exploring the large space of complex behaviors. We have previously shown that this approach can be effective for certains kinds of systems in biotechnology [2], [15] and medicine [13]. Our challenge is to adapt incremental, hierarchical refinement to modeling organisms and communities in metagenomic and comparative genomic applications.

# 4. Application Domains

## 4.1. Genome and transcriptome annotation, to model function

Sequencing genomes and transcriptomes provides a picture of how a biological system can function, or does function under a given physiological condition. Simultaneous sequencing of a group of related organisms is now a routine procedure in biological laboratories for studying a behavior of interest, and provides a marvelous opportunity for building a comprehensive knowledge base of the relations between genomes. Key elements in mining these relations are: classifying the genes in related organisms and the reactions in their metabolic networks, recognizing the patterns that describe shared features, and highlighting specific differences.

PLEIADE will develops applications in comparative genomics of related organisms, using new mathematical tools for representing compactly, at different scales of difference, comparisons between related genomes. New methods based on computational geometry refine these comparisons. Compact representations can be stored, exchanged, and combined. They will form the basis of new simultaneous genome annotation methods, linked directly to abductive inference methods for building functional models of the organisms and their communities.

Our ambition in biotechnology is to permit the design of synthetic or genetically selected organisms at an abstract level, and guide the modification or assembly of a new genome. Our effort is focused on two main applications: genetic engineering and synthetic biology of oil-producing organisms (biofuels in CAER, palm oils), and improving and selecting starter microorganisms used in winemaking (collaboration with the ISVV and the BioLaffort company).

### 4.1.1. *Oil Palm lipid synthesis*

The largest source of vegetable oil [1] is the fruit mesocarp of the oil palm *Elaeis guineensis*, a remarkable tissue that can accumulate up to 90% oil, the highest level observed in the plant kingdom. The market share of oil palm is expected to increase in order to meet increased demand for vegetable oil, predicted to double by 2030 [18], be it as food or as a source of biofuels in Africa. A significant proportion of palm oil is produced on small estates that do not have access to efficient milling facilities, and run a great risk of spoilage through oil acidification. Improving palm oil quality through genetics and selection will result in economic gains [24] by addressing several targets such as improvement of oil yield, tuning of oil quality through the rate of unsaturated fatty acids or impairment of degradation processes. Furthermore, as genome biodiversity resides mostly in Africa, oil from African oil palms can vary greatly in fatty acid composition according to cultivar genetic differences and to weather conditions, and the precise mechanisms regulating this variability are not yet understood.

---

[1] 32% of the world market share [24]

A growing body of molecular resources for studying oil palm fruit are making it possible to study and improve the quality and quantity of oil produced by oil palms. In particular, these oils can vary greatly in fatty acid composition, and while the precise mechanisms regulating this variability are not completely understood, establishing a link between oil palm genotype and phenotype appears increasingly feasible. PLEIADE will work with the CNRS/UB UMR 5200 (LBM), a laboratory with an established reputation in studying fatty acid metabolism in *E. guineensis*, to improve understanding of the links between genetic diversity and oil production, and participate in developing applications.

### 4.1.2. Engineering pico-algae



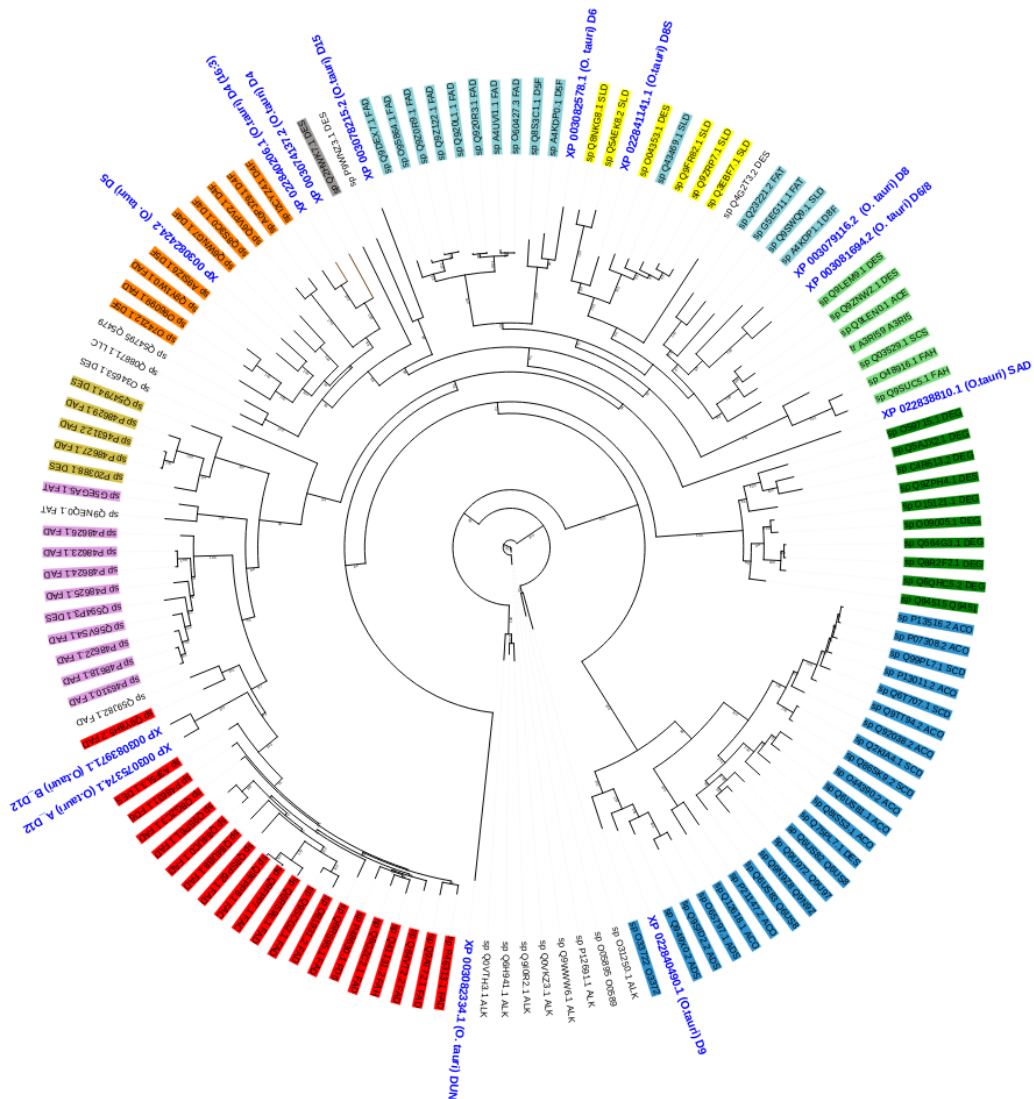Figure 3. Phylogenetic structure of long-chain polyunsaturated fatty acid desaturase specificity. Highlighted are thirteen desaturases from Ostreococcus tauri

Docosahexaenoic acid (DHA) is an essential nutriment for human brain tissue and can only be obtained from marine or riverine fish that live on phytoplankton and zooplankton, since human neurons lack the delta desaturase required for *de novo* synthesis of DHA. [21] Unfortunately, fishing is become less and less a sustainable resource. Since phytoplankton and zooplankton are the ultimate source of DHA consumed, there is considerable interest in obtaining DHA directly rather than through the intermediary of fish. A very promising approach is through the bio-engineering of pico-algae.

In order to produce the long-chain polyunsaturated fatty acids (LC-PUFA) needed for human nutrition, it is necessary to precisely engineer the desaturases that produce them. Desaturases are enzymes responsible for the introduction of double bonds into fatty acids. Desaturases are specific in recognizing their substrates and in placing the double bond in the proper place. The desaturases that produce the LC-PUFA necessary for human nutrition are present only in some species.

Our goal is to design methods to predict the substrate and region specificities for desaturases in algal species, particularly *Ostreococcus tauri*, the smallest photosynthetic eukaryote that can be cultivated. Thirteen desaturases are known in *O. tauri* and can be placed in the phylogeny of the desaturase family (figure 3. The biochemical and structural characterization of these enzymes is as yet very incomplete. This work is ongoing (see section 8.2.3) and requires close collaboration between biologists and computer scientists.

## 4.2. Molecular based systematics and taxonomy

Defining and recognizing myriads of species in biosphere has taken phenomenal energy over the past centuries and remains a major goal of Natural History. It is an iconic paradigm in pattern recognition (clustering has coevolved with numerical taxonomy many decades ago). Developments in evolution and molecular biology, as well as in data analysis, have over the past decades enabled a profound revolution, where species can be delimited and recognized by data analysis of sequences. We aim at proposing new tools, in the framework of E-science, which make possible (*i*) better exploration of the diversity in a given clade, and (*ii*) assignment of a place in these patterns for new, unknown organisms, using information provided by sets of sequences. This will require investment in data analysis, machine learning, and pattern recognition to deal with the volumes of data and their complexity.

One example of this project is about the diversity of trees in Amazonian forest, in collaboration with botanists in French Guiana. Protists (unicellular Eukaryots) are by far more diverse than plants, and far less known. Molecular exploration of Eukaryotes diversity is nowadays a standard in biodiversity studies. Data are available, through metagenomics, as an avalanche and make molecular diversity enter the domain of Big Data. Hence, an effort will be invested, in collaboration with other Inria teams (GenScale, HiePACS) for porting to HPC algorithms of pattern recognition and machine learning, or distance geometry, for these tools to be available as well in metagenomics. This will be developed first on diatoms (unicellular algae) in collaboration with INRA team at Thonon and University of Uppsala), on pathogens of tomato and grapewine, within an existing network, and on bacterial communities, in collaboration with University of Pau. For the latter, the studies will extend to correlations between molecular diversity and sets of traits and functions in the ecosystem.

## 4.3. Community ecology and population genetics

Community assembly models how species can assemble or diassemble to build stable or metastable communities. It has grown out of inventories of countable organisms. Using *metagenomics* one can produce molecular based inventories at rates never reached before. Most communities can be understood as pathways of carbon exchange, mostly in the form of sugar, between species. Even a plant cannot exist without carbon exchange with its rhizosphere. Two main routes for carbon exchange have been recognized: predation and parasitism. In predation, interactions–even if sometimes dramatic–may be loose and infrequent, whereas parasitism requires what Claude Combes has called intimate and sustainable interactions [19]. About one decade ago, some works [25] have proposed a comprehensive framework to link the studies of biodiversity with community assembly. This is still incipient research, connecting community ecology and biogeography.

We aim at developing graph-based models of co-occurence between species from NGS inventories in metagenomics, i.e. recognition of patterns in community assembly, and as a further layer to study links, if any, between diversity at different scales and community assemblies, starting from current, but oversimplified theories, where species assemble from a regional pool either randomly, as in neutral models, or by environmental filtering, as in niche modeling. We propose to study community assembly as a multiscale process between nested pools, both in tree communities in Amazonia, and diatom communities in freshwaters. This will be a step towards community genomics, which adds an ecological flavour to metagenomics.

Convergence between the processes that shape genetic diversity and community diversity–drift, selection, mutation/speciation and migration–has been noted for decades and is now a paradigm, establishing a continuous scale between levels of diversity patterns, beyond classical approaches based on iconic levels like species and populations. We will aim at deciphering diversity pattern along these gradients, connecting population and community genetics. Therefore, some key points must be adressed on reliability of tools.

Next-generation sequencing technologies are now an essential tool in population and community genomics, either for making evolutionary inferences or for developing SNPs for population genotyping analyses. Two problems are highlighted in the literature related to the use of those technologies for population genomics: variable sequence coverage and higher sequencing error in comparison to the Sanger sequencing technology. Methods are developed to develop unbiased estimates of key parameters, especially integrating sequencing errors [23]. An additional problem can be created when sequences are mapped on a reference sequence, either the sequenced species or an heterologous one, since paralogous genes are then considered to be the same physical position, creating a false signal of diversity [20]. Several approaches were proposed to correct for paralogy, either by working directly on the sequences issued from mapped reads [20] or by filtering detected SNPs. Finally, an increasingly popular method (RADseq) is used to develop SNP markers, but it was shown that using RADseq data to estimate diversity directly biases estimates [12]. Workflows to implement statistical methods that correct for diversity biases estimates now need an implementation for biologists.

# 5. Highlights of the Year

## 5.1. Highlights of the Year

Pascal Durrens was promoted this year to the new grade of *Chargé de Recherche Hors Classe* of the CNRS.

# 6. New Software and Platforms

## 6.1. Magus

KEYWORDS: Bioinformatics - Genomic sequence - Knowledge database
SCIENTIFIC DESCRIPTION: MAGUS can be used on small installations with a web server and a relational database on a single machine, or scaled out in clusters or elastic clouds using Apache Cassandra for NoSQL data storage and Apache Hadoop for Map-Reduce.
FUNCTIONAL DESCRIPTION: The MAGUS genome annotation system integrates genome sequences and sequences features, in silico analyses, and views of external data resources into a familiar user interface requiring only a Web navigator. MAGUS implements annotation workflows and enforces curation standards to guarantee consistency and integrity. As a novel feature the system provides a workflow for simultaneous annotation of related genomes through the use of protein families identified by in silico analyses this has resulted in a three-fold increase in curation speed, compared to one-at-a-time curation of individual genes. This allows us to maintain standards of high-quality manual annotation while efficiently using the time of volunteer curators.

NEWS OF THE YEAR: Magus is now available as a Docker image, and can be integrated with other containerized services using Pleiade's Alcyone system.

- Participants: David Sherman, Florian Lajus, Natalia Golenetskaya, Pascal Durrens and Xavier Calcas
- Partners: Université de Bordeaux - CNRS - INRA
- Contact: David James Sherman
- Publication: High-performance comparative annotation
- URL: http://magus.gforge.inria.fr

## 6.2. Mimoza

KEYWORDS: Systems Biology - Bioinformatics - Biotechnology

FUNCTIONAL DESCRIPTION: Mimoza uses metabolic model generalization and cartographic paradigms to allow human experts to explore a metabolic model in a hierarchical manner. Mimoza generalizes genome-scale metabolic models, by factoring equivalent reactions and metabolites while preserving reaction consistency. The software creates an zoomable representation of a model submitted by the user in SBML format. The most general view represents the compartments of the model, the next view shows the visualization of generalized versions of reactions and metabolites in each compartment , and the most detailed view visualizes the initial model with the generalization-based layout (where similar metabolites and reactions are placed next to each other). The resulting map can be explored on-line, or downloaded in a COMBINE archive. The zoomable representation is implemented using the Leaflet JavaScript library for mobile-friendly interactive maps. Users can click on reactions and compounds to see the information about their annotations.

NEWS OF THE YEAR: Mimoza is now available as a Docker image, and can be integrated with other containerized services using Pleiade's Alcyone system.

- Participants: Anna Zhukova and David James Sherman
- Contact: David James Sherman
- Publications: Knowledge-based generalization of metabolic models - Knowledge-based zooming for metabolic models - Knowledge-based generalization of metabolic networks: a practical study
- URL: http://mimoza.bordeaux.inria.fr/

## 6.3. Declic

FUNCTIONAL DESCRIPTION: Declic is a Python library that provides several tools for data analysis in the domains of multivariate data analysis, machine learning, and graph based methods. It can be used to study in-depth the accuracy of the dictionary between molecular based and morphological based taxonomy.

Declic includes an interpreter for a Domain Specific Language (DSL) to make its Python library easy to use for scientists familiar with environments such as R.

- Partner: INRA
- Contact: Alain Franc

## 6.4. Diagno-Syst

*diagno-syst: a tool for accurate inventories in metabarcoding*

KEYWORDS: Biodiversity - Clustering - Ecology

FUNCTIONAL DESCRIPTION: Diagno-syst builds accurate inventories for biodiversity. It performs supervised clustering of reads obtained from a next-generation sequencing experiment, mapping onto an existing reference database, and assignment of taxonomic annotations.

- Participants: Alain Franc, Jean-Marc Frigerio, Philippe Chaumeil and Franck Salin
- Partner: INRA
- Contact: Alain Franc
- Publication: diagno-syst: a tool for accurate inventories in metabarcoding

## 6.5. Alcyone

*Alcyone instantiates bioinformatics environments from specifications committed to a Git repository*

KEYWORDS: Docker - Orchestration - Bioinformatics - Microservices - Versioning

SCIENTIFIC DESCRIPTION: Alcyone conceives the user's computing environment as a microservices architecture, where each bioinformatics tool in the specification is a separate containerized Docker service. Alcyone builds a master container for the specified environment that is responsible for building, updating, deploying and stopping these containers, as well as recording and sharing the environment in a Git repository. The master container can be manipulated using a command-line interface.

FUNCTIONAL DESCRIPTION: Alcyone defines a file structure for the specifying bioinformatics analysis environments, including tool choice, interoperability, and sources of raw data. These specifications are recorded in a Git repository. Alcyone compiles a specification into a master Docker container that deploys and orchestrates containers for each of the component tools. Alcyone can restore any version of an environment recorded in the Git repository.

NEWS OF THE YEAR: Alcyone was designed and implemented this year.

- Participants: Louise-Amelie Schmitt and David Sherman
- Contact: David Sherman
- URL: https://team.inria.fr/pleiade/alcyone/

## 6.6. family-3d

KEYWORDS: Biodiversity - Point cloud - 3D modeling

SCIENTIFIC DESCRIPTION: The method statistically selects a subset of pairwise distances between proteins in the family, constructs a weighted graph, and lays it out using an adaptation of the three-dimensional extension of the Kamada-Kawai force-directed layout.

FUNCTIONAL DESCRIPTION: Family-3D lays out high-dimension protein family point clouds in 3D space. The resulting lower-dimension forms can be printed, so that they can be explored and compared manually. They can also be explored interactively or stereographically.

Comparison of the 3D forms reveals classes of structurally similar families, whose characteristic shapes correspond to different evolutionary scenarios. Some of these scenarios are: neofunctionalization, subfunctionalization, founder gene effect, ancestral family.

To facilitate curator training, Family-3D includes an interactive terminal containing a microcontroller, an RFID reader, and an LED ring. A set of shapes that fall in predetermined classes is printed, with a unique RFID tag in each shape. Trainees classify family shapes by manual inspection and submit their classes to the terminal, which evaluates the proposed class and provides visual feedback.

- Participant: David Sherman
- Contact: David Sherman
- URL: https://gitlab.inria.fr/pleiade/family-3d

## 6.7. magecal

KEYWORD: Genomics

SCIENTIFIC DESCRIPTION: Magecal independently runs training and prediction steps for Augustus, Conrad, GeneID, GeneMark, and Snap. The results are cleaned and integrated into a common format. Jigsaw is trained and used for model reconciliation. Consistency constraints are applied to ensure that phase and intron structure are biologically plausible.

FUNCTIONAL DESCRIPTION: Magecal predicts a set of protein coding genes in fungal genomic sequences, using different de novo prediction algorithms, and reconciling the predictions with the aid of comparative data. Magecal applies consistency constraints to guarantee that the predicted genes are biologically valid.

RELEASE FUNCTIONAL DESCRIPTION: Dockerization and compatibility with Alcyone

- Participants: Pascal Durrens and David Sherman
- Contact: David Sherman
- URL: https://gitlab.inria.fr/magecal/magecal

## 6.8. AsebaHub

*Turn-key bridging of Aseba mobile robots to wifi networks*
KEYWORD: Robotics
SCIENTIFIC DESCRIPTION: Aseba is an open-source modular architecture for event-based control of multi-microcontroller robots, used in the Thymio-II educational robot, and integrated into development environments for teaching robotics programming to 8-18 year old children. Since 2014 Inria has contributed several modules to Aseba, including mDNS-sd support for discovering TCP/IP connections, an HTTP bridge that allows Aseba robots to be used with the Scratch programming environment, and photo-realistic 2D simulation.

Thymio-II robots only communicate through USB connections from a host computer, so until now it has been necessary to install Aseba software on that computer. This is not always possible in schools.

AsebaHub is a firmware image for wifi routers that takes responsibility for USB communication with Aseba robots. Each robot is exposed as a TCP/IP target, advertised using mDNS-sd. AsebaHub can group robots together in ad hoc Aseba networks. It provides HTTP access to robots so that they can be used with Scratch or Snap!. These services work out of the box, without any host configuration or software installation.

AsebaHub can operate as a wifi Access Point, or as a network Bridge to an existing local-area network. A slide switch on the side of the router seamlessly switches between these configurations.

AsebaHub is implemented by compiling Aseba software for the OpenWrt/LEDE open source router platform, adding interfacing code that automatically associates new robots, starts and stops the corresponding network services, and provides a web-based administrative interface.
FUNCTIONAL DESCRIPTION: AsebaHub is a small device with a USB port for connecting Aseba robots, single or with a USB hub. Every robot is made available as a network target on the wifi and wired local-area networks, where they can be discovered using mDNS-sd (Zeroconf/Bonjour). No network or user configuration is necessary.

AsebaHub works out of the box with the Aseba Studio and VPL programming environments. AsebaHub can be configured to also provide an HTTP service that is compatible with Scratch or other environments.

AsebaHub acts as either a wireless Access Point to which users may connect, or as a Bridge to an existing local-area network. Either configuration can be chosen seamlessly with the flick of a switch.

- Participants: David Sherman and Louise-Amelie Schmitt
- Contact: David Sherman

# 7. New Results

## 7.1. Introgressions as a source of diversity

Several prominent mechanisms of genomic evolution have been described for the yeasts, among them inter-specific hybridization, reticulated evolution, aneuploidization, recent or ancient poly-ploidization events, large chromosomal duplication or more limited gene duplication, and horizontal transfer. These mechanisms are usually so closely intertwined that it is difficult to determine which ones are causes or consequences. Regardless of mechanisms the result has been a drastic reshaping of yeasts genome along evolution. Understanding these mechanisms is important, not only for strain construction in biotechnology, but also more fundamentally for insight into the causes and effects of genome reshaping on much shorter time scales.

Introgression, the transfer of large or more limited genetic information from one species to another, is an evolutionary mechanism of particular interest in industrial applications such as wine making where large vat cultures are used. Introgression results in mosaic genomes, and can be the result of interspecific hybridization fol- lowed by the extensive loss of one parental genome, either through repeated backcross with one parental species or through missegregation of the hybrid at meiosis.

In collaboration with the Institut des Science de la Vigne et du Vin and Bordeaux Sciences Agro, PLEIADE developed tools to rapidly assess the presence of introgressed regions in a large population of *Saccharomyces uvarum* isolates (104 strains), focusing on Holarctic isolates from natural, cider and wine environments since introgressed regions are absent in Southern hemisphere isolates. The overall number of introgressed regions is significantly higher in cider-associated strains compared to wild strains, and is higher in wine isolates. However, only a subset of the introgressed regions were found to be overrepresented in anthropic activities and their number and quality varied between cider- and wine-making processes.

Paradoxically, the low Holarctic genetic diversity observed in [1] contrasts with the relative high phenotypic diversity found for technological traits. This contradiction suggests that interspecific introgressions found among Holarctic *S. uvarum* strains could be the most important source of genetic, and by extention of phenotypic, diversity.

## 7.2. Geometric view of biodiversity and molecular taxonomy

For the geometric view on biodiversity and molecular based taxonomy (metabarcoding), 2018 has been characterized by significant progress in the connections between the questions posed by metabarcoding, and the implementaton of solutions by intensive computing in collaboration with HiePACS team. We are now able to run a Multidimensional Scaling (MDS) to build a point cloud where each point is a read (a short sequence from an environmental sample) on a sample of nearly $2 \times 10^5$ reads on one node of a cluster with shared memory. PLEIADE is member of two projects lead by Inria BSO to extend these possibilities to work with distributed memory (connection with the Chameleon dense linear algebra solver) with the possibility to reach results with one million sequences at hand. To our knowledge, such a connection between scientific questions related to biodiversity and HPC-HTC has up to now had no equivalent and is new.

In collaboration with HiePACS, the SED of Inria BSO, and the GRICAD (Mésocentre of Grenoble), PLEIADE has developed a scientific library of tools for dimension reductions in the framework of Big Data, specifically linear algebra of full matrices. A prototype was developed in numpy, then reimplemented in C++ with an external library for Random Projection, and made available in Julia. This library will be connected with the workflow developed in Regional project and ADT Gordon.

## 7.3. Tangible exploration of protein families

Protein families[3] are an effective way to compare the complete genomes of fungal species. In general these comparisons are very challanging due to the large evolutionary distances involved, the wide range of GC compositions observed from one species to the next, and the extensive map reshuffling that is characteristic of the yeasts in particular. Protein families are a classification of protein-coding gene sequences into phylogenetic groups, using clustering methods and semi-supervised classification. Members of a family are homologous and in many cases this homology is suggestive of functional similarity.

An intriguing feature of protein families is that the weighted graph constructed from their pairwise distance matrices has a structure that reflects the evolutionary history of the family. We developed software (family-3d, see above) that uses truncated distances to construct a weighted graph and to lay it out using an adaptation of the three-dimensional extension of the Kamada-Kawai force-directed layout. The resulting shapes for a set of protein families are then clustered manually by similarity (figure 5). Similarity in shape is highly suggestive of similarity in evolutionary scenarios, leading to hypotheses about the histories of individual protein families and the mechanisms by which functional diversity is obtained.
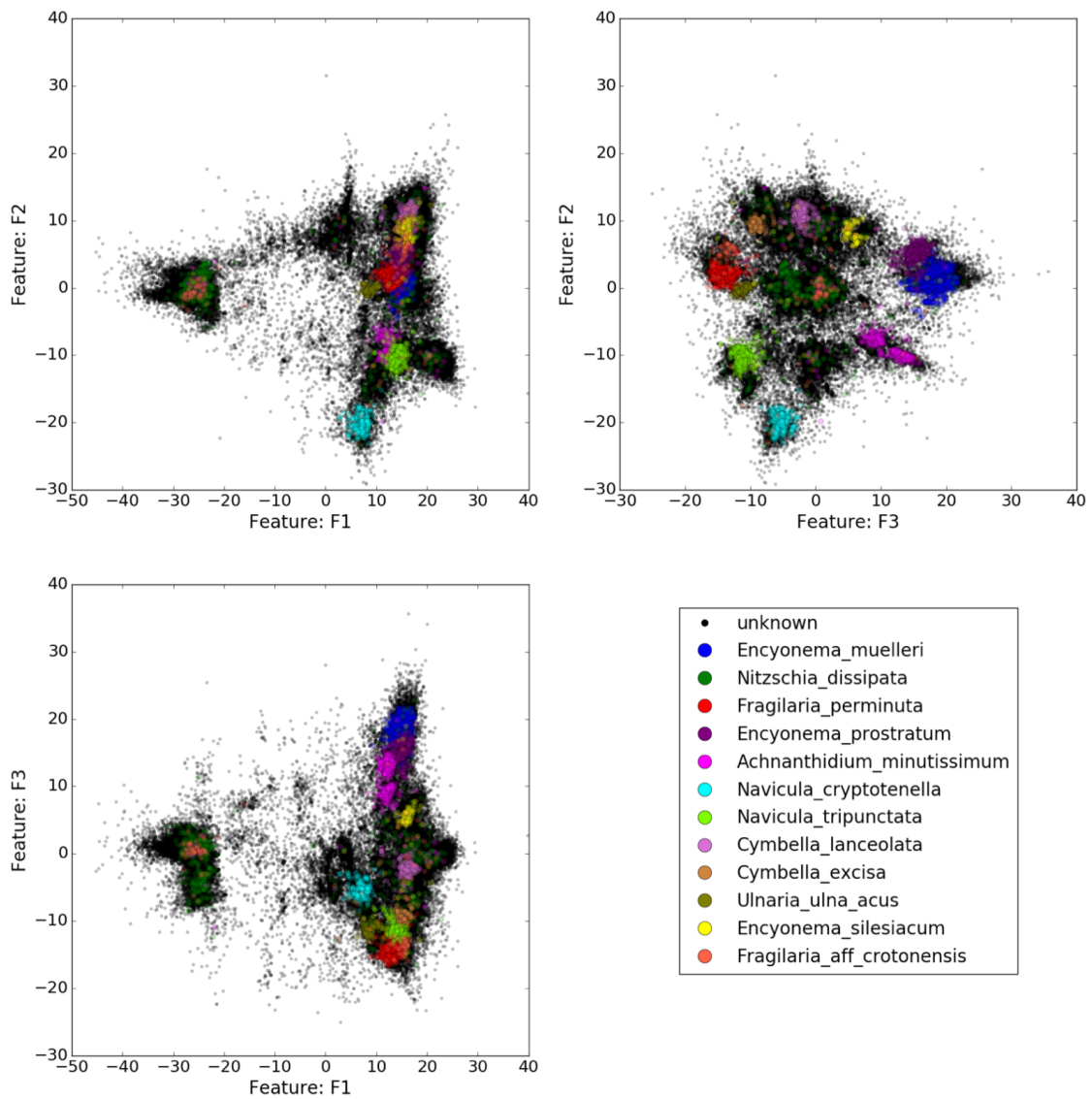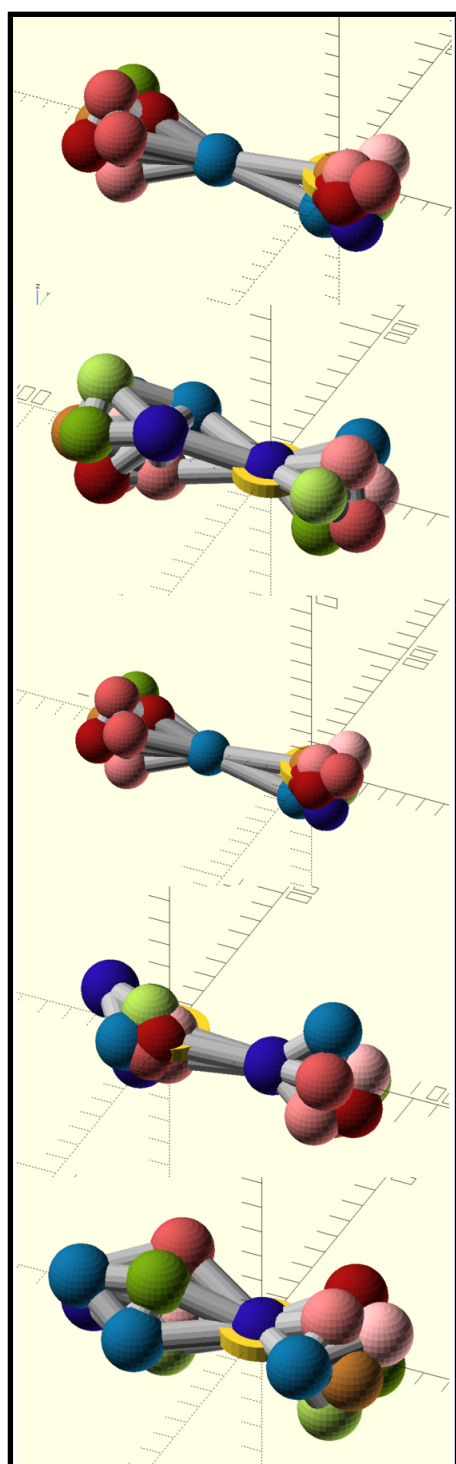
*Figure 4. Validation of high density islands using supervised classification. Metagenomic reads from diatoms in Lake Geneva [26] were analyzed by the method from [7] and colored by species according to a reference database.*

a) subfunctionalization

b) articulation

c) neofunctionalization

*Figure 5. Three examples of classes of protein families with similar shapes: a) articulation, b) subfunctionalization, c) neofunctionalization*

# 8. Partnerships and Cooperations

## 8.1. Regional Initiatives

### 8.1.1. High-performance computing and metabarcoding

PLEIADE is member of two projects, one funded by the Région Nouvelle Aquitaine and one funded as Inria ADT Gordon, connecting Chameleon, StartPU and NewMadeleine, where the use case of metabarcoding (questions, data sets) hase been selected to link these layers together. This will permit us to address unsupervised clustering of one million reads next year. These projects are in collaboration with the HiePACS, TADAAM, and STORM project-teams.

### 8.1.2. COTE – Continental to Coastal Ecosystems

The Labex cluster of excellence COTE (Continental To coastal Ecosystems: evolution, adaptability and governance) develops tools to understand and predict ecosystem responses to human-induced changes as well as methods of adaptative management and governance to ensure their sustainability. The LabEx includes nine laboratories of the University of Bordeaux and major national research institutes involved in research on terrestrial and aquatic ecosystems (INRA, CNRS, IFREMER and IRSTEA). PLEIADE is a partner in two projects funded by COTE:

- *Malabar: Métabarcoding appliqué à la biodiversté des communauéts algales du Bassin d'Arcachon*

- *Aerobarcoding: détection de pollens allergenisants*. 2017-18.

## 8.2. National Initiatives

### 8.2.1. Biocontamination in aircraft reservoirs

ANTICOR is an industrial-academic research and development working group coordinated by Dassault Aviation, investigating the causes of microbial contamination in aircraft reservoirs and aimed at developing mitigating procedures and equipment. Previous results have shown that this contamination forms biofilms at the fuel-water interface and is comprised of complex communities of hundreds of bacterial and fungal species. PLEIADE is particularly interested in measuring and modeling these communities, especially as concerns understanding how they change based on environmental conditions and on reservoir geometry.

This working group continues work started in CAER – Alternative Fuels for Aeronautics, a 6 M-Euro contract with the Civil Aviation Directorate (Direction Générale de l'Aviation Civile, DGAC), coordinated by the French Petroleum Institute (Institut français de pétrole-énergies nouvelles, IFPEN) on behalf of a large consortium of industrial (EADS, Dassault, Snecma, Turbomeca, Airbus, Air France, Total) and academic (CNRS, INRA, Inria) partners to explore different technologies for alternative fuels for aviation.

### 8.2.2. Agence Française pour la Biodiversité

The AFB is a public law agency of the French Ministry of Ecology that supports public policy in the domains of knowledge, preservation, management, and restoration of biodiversity in terrestrial, aquatic, and marine environments. PLEIADE is a partner in two AFB projects developed with the former ONEMA:

- *Methods for metabarcoding*. 2017-18.
- *Molecular diagnosis of freshwater quality*. 2014-present.

### 8.2.3. Inria Projet Lab in silico Algae

In 2017 PLEIADE joined the IPL "In silico Algae" coordinated by Olivier Bernard. The IPL addresses challenges in modeling and optimizing microalgae growth for industrial applications. PLEIADE worked this year on comparative genomic analysis of genes implicated in lipid production by the picoalgae *Ostreococcus tauri*, in collaboration with Florence Corellou of the CNRS UMR 5200 (Laboratoire de Biogénèse Membranaire). The goal of this work is the production of long-chain polyunsaturated fatty acids, developed as nutritional additives. Mercia Ngoma Komb's two-month internship in PLEIADE contributed to this work.

## 8.3. European Initiatives

### 8.3.1. Collaborations in European Programs, Except FP7 & H2020

Alain Franc is co-chair of Working Group 4 (Data Analysis and Storage) of COST DNAqua.net [2], with the main task of developing contact with HPC and metabarcoding for serving the whole community. He chaired the meeting of this groups in Pécs in 2018. The goal of DNAqua-Net is to nucleate a group of researchers across disciplines with the task to identify gold-standard genomic tools and novel eco-genomic indices and metrics for routine application for biodiversity assessments and biomonitoring of European water bodies.

## 8.4. International Initiatives

### 8.4.1. Vitapalm – Food and nutrition security and sustainable agriculture in Africa

PLEIADE has joined the Vitapalm program financed by LEAP-Agri [3], the joint Europe Africa Research and Innovation (R&I) initiative related to Food and Nutrition Security and Sustainable Agriculture. Vitapalm will use genomics and selection to improve the nutritional quality and the stability of palm oil produced by Africa smallholdings for local consumption. Project partners are from Cameroon, France, Germany, and Ghana.

### 8.4.2. Simulation of metacommunities

In collaboration with the Pasteur Institute in Cayenne and the INRA MIA Research Team in Toulouse, PLEIADE is developing a stochastic model for simulation of metacommunities, in the framework of patch occupancy models. The objective is a better understanding of zoonose propagation, namely rabies through bat hosts in connection with disturbances of pristine forests in French Guiana, which have an impact on the exposure of human populations to wildlife that act as reservoirs of zoonoses.

### 8.4.3. CEBA – Center for the study of biodiversity in Amazonia

The Laboratoire of excellence CEBA promotes innovation in research on tropical biodiversity. It brings together a network of internationally-recognized French research teams, contributes to university education, and encourages scientific collaboration with South American countries. PLEIADE participates in three current international projects funded by CEBA:

- *MicroBIOMES: Microbial Biodiversities*. 2017-19.
- *Neutrophyl: Inferring the drivers of Neotropical diversification*. 2017-19.
- *Phyloguianas: Biogeography and pace of diversification in the Guiana Shield*. 2015-present

# 9. Dissemination

## 9.1. Promoting Scientific Activities

### 9.1.1. Scientific Events Selection

#### 9.1.1.1. Member of the Conference Program Committees

Alain Franc co-organized with IRSTEA Bordeaux, IFREMER Nantes and INRA Toulouse the biannual meeting of the Mexico Network at Inria, Ada Lovelace room, in November 2018 (see https://mexico2018.sciencesconf.org/) and s a member of the Scientific Committee. The theme of the meeting was "Optimization, Sensitivity Analysis, and Exploration of Model Outputs." Mexico is a methodological network dedicated to methods of computer exploration of complex models. [4]

---

[2]http://dnaqua.net/
[3]http://www.leap-agri.com/
[4]http://reseau-mexico.fr

### *9.1.2. Journal*

*9.1.2.1. Member of the Editorial Boards*

Alain Franc is member of the editorial board of BMC Evolutionary Biology.

Pascal Durrens is a member of the editorial board of the journal ISRN Computational Biology.

### *9.1.3. Research Administration*

David Sherman is a co-president of the INRA Commission for the Evaluation of Engineers (CEI IMP6) [5]. Every INRA staff engineer is evaluated by this commission every four years.

Alain Franc has been appointed "chargé de mission calcul" at INRA by INRA Delegate for Digital Transition. As such, his mission is to propose animations and solutions for the development of scientific computing at INRA, whatever the Research Department.

## 9.2. Popularization

### *9.2.1. Internal or external Inria responsibilities*

David Sherman is president of the Commission for Technology Development (CDT) of the Inria Bordeaux Sud-Ouest research center. The CDT has two roles. First, it evaluates funding requests for Technology Development and Technology Transfer projects, which typically involve hiring technical staff. Second, the CDT is responsible for validating and overseeing contract engineers hired by Inria project-teams.

### *9.2.2. Interventions*

The "Family-3D" activity (see below) was presented on three occasions:
- During the 10 year celebration on the evening of September 27
- During the Open House all day October 13
- During the Fête de la Science, to 8 groups of 12 students

### *9.2.3. Creation of media or tools for science outreach*

Pleiade created an interactive activity "Family-3D" to explain the methods and uses of pattern classification of protein families. Using new software developed by the team we laid out 30 families in 3D space and converting those shapes to 3D volumes. The selected families had been previously grouped into biologically pertinent classes by human curators. Each shape was printed with an embedded unique RFID tag. We also designed and printed five interactive terminals, containing a microcontroller, an RFID reader with a custom-designed inductive coil as input, and an LED ring as output. Participants would propose groups of shapes that they believed belonged to the same class, and the terminal would evaluate the group.

David Sherman contributes open-source software development to the Aseba platform for educational robotics [6], deployed in Thymio II robots used by children as well as in the simulator used by Class'Code [7] to train teachers.

David Sherman and Louise-Amélie Schmitt developed AsebaHub, a specialized WiFi hub for connecting Thymio-II robots into a local-area network. AsebaHub acts as either a wireless Access Point to which users may connect, or as a Bridge to an existing local-area network. Robots can be discovered using mDNS-sd (Zeroconf/Bonjour). AsebaHub works out of the box with the Aseba Studio and VPL programming environments, and can also provide an HTTP service that is compatible with Scratch, Snap!, or other environments.

---

[5] INRA UAR0837 DEV Délégation à l'Evaluation
[6] http://aseba.io/
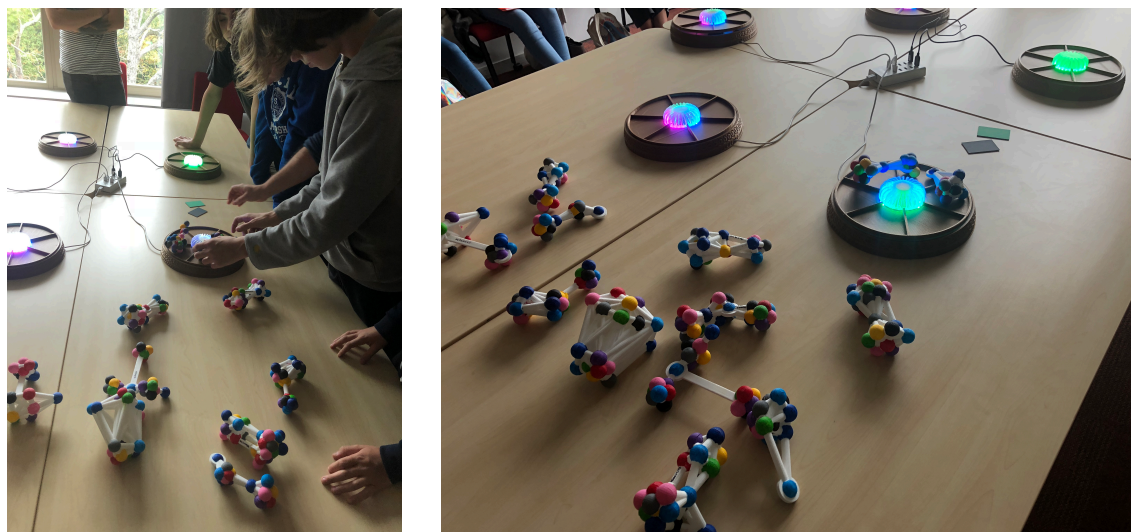[7] https://pixees.fr/classcode-la-formation-associee-a-pixees/

*Figure 6. Family-3D activity during the Circuit Scientifique of the Fête de la Science, 2018*

# 10. Bibliography

## Major publications by the team in recent years

[1] P. ALMEIDA, C. GONÇALVES, S. TEIXEIRA, D. LIBKIND, M. BONTRAGER, I. MASNEUF-POMAREDE, W. ALBERTIN, P. DURRENS, D. J. SHERMAN, P. MARULLO, C. TODD HITTINGER, P. GONÇALVES, J. P. SAMPAIO. *A Gondwanan imprint on global diversity and domestication of wine and cider yeast Saccharomyces uvarum*, in "Nature Communications", 2014, vol. 5, 4044 p. [*DOI :* 10.1038/NCOMMS5044], https://hal.inria.fr/hal-01002466

[2] R. ASSAR, M. A. MONTECINO, A. MAASS, D. J. SHERMAN. *Modeling acclimatization by hybrid systems: Condition changes alter biological system behavior models*, in "BioSystems", June 2014, vol. 121, pp. 43-53 [*DOI :* 10.1016/J.BIOSYSTEMS.2014.05.007], https://hal.inria.fr/hal-01002987

[3] D. J. SHERMAN, T. MARTIN, M. NIKOLSKI, C. CAYLA, J.-L. SOUCIET, P. DURRENS. *Genolevures: protein families and synteny among complete hemiascomycetous yeast proteomes and genomes*, in "Nucleic Acids Research", 2009, vol. 37, pp. D550-D554 [*DOI :* 10.1093/NAR/GKN859], https://hal.inria.fr/inria-00341578

## Publications of the year

### Articles in International Peer-Reviewed Journals

[4] W. ALBERTIN, M. CHERNOVA, P. DURRENS, E. GUICHOUX, D. J. SHERMAN, I. MASNEUF-POMAREDE, P. MARULLO. *Many interspecific chromosomal introgressions are highly prevalent in Holarctic Saccharomyces*

*uvarum strains found in human-related fermentations*, in "Yeast", January 2018, vol. 35, n[o] 1, pp. 141-156 [*DOI :* 10.1002/YEA.3248], https://hal.inria.fr/hal-01585808

[5] F. KECK, A. A. FRANC, M. KAHLERT. *Disentangling the processes driving the biogeography of freshwater diatoms: A multiscale approach*, in "Journal of Biogeography", July 2018, vol. 45, n[o] 7, pp. 1582-1592, https://hal.inria.fr/hal-01982135

### Scientific Books (or Scientific Book chapters)

[6] F. LEESE, A. BOUCHEZ, K. ABARENKOV, F. ALTERMATT, A. BORJA, K. BRUCE, T. EKREM, F. ČIAMPOR, Z. ČIAMPOR, F. COSTA, S. DUARTE, V. ELBRECHT, D. FONTANETO, A. A. FRANC, M. GEIGER, D. HERING, M. KAHLERT, B. KALAMUJIĆ STROIL, M. KELLY, E. KESKIN, I. LISKA, P. MERGEN, K. MEISSNER, J. PAWLOWSKI, L. PENEV, Y. REYJOL, A. ROTTER, D. STEINKE, B. VAN DER WAL, S. S. VITECEK, J. ZIMMERMANN, A. WEIGAND. *Why We Need Sustainable Networks Bridging Countries, Disciplines, Cultures and Generations for Aquatic Biomonitoring 2.0: A Perspective Derived From the DNAqua-Net COST Action*, in "Next Generation Biomonitoring: Part 1", Elsevier, 2018, vol. 58, pp. 63-99, https://hal.inria.fr/hal-01984996

### Research Reports

[7] P. BLANCHARD, P. CHAUMEIL, J.-M. FRIGERIO, F. RIMET, F. SALIN, S. THÉROND, O. COULAUD, A. A. FRANC. *A geometric view of Biodiversity: scaling to metagenomics*, Inria ; INRA, January 2018, n[o] RR-9144, pp. 1-16, https://arxiv.org/abs/1803.02272 , https://hal.inria.fr/hal-01685711

[8] A. A. FRANC, P. BLANCHARD, O. COULAUD. *Nonlinear Mapping and Distance Geometry*, Inria Bordeaux Sud-Ouest, September 2018, n[o] RR-9210, 14 p. , https://arxiv.org/abs/1810.08661 , https://hal.inria.fr/hal-01897104

[9] R. ISSA, D. J. SHERMAN. *Deriving Gene Associations for Metabolic Models*, Inria Bordeaux Sud-Ouest, December 2018, n[o] RR-9246, https://hal.inria.fr/hal-01976744

### Other Publications

[10] T. LANG, P. P. ABADIE, V. LÉGER, T. DECOURCELLE, J.-M. FRIGERIO, C. BURBAN, C. BODENES, E. GUICHOUX, G. G. LE PROVOST, C. ROBIN, N. TANI, P. P. LÉGER, C. LEPOITTEVIN, V. A. EL MUJTAR, F. HUBERT, J. TIBBITS, J. PAIVA, A. A. FRANC, F. RASPAIL, S. MARIETTE, M.-P. M. REVIRON, C. PLOMION, A. KREMER, M.-L. DESPREZ-LOUSTAU, P. GARNIER-GERE. *High-quality SNPs from genic regions highlight introgression patterns among European white oaks (Quercus petraeaandQ. robur)*, August 2018, working paper or preprint [*DOI :* 10.1101/388447], https://hal.inria.fr/hal-01985021

## References in notes

[11] R. ALUR. *SIGPLAN Notices*, in "Generating Embedded Software from Hierarchical Hybrid Models", 2003, vol. 38, n[o] 7, pp. 171–82

[12] B. ARNOLD, R. CORBETT-DETIG, D. HARTL, K. BOMBLIES. *RADseq underestimates diversity and introduces genealogical biases due to nonrandom haplotype sampling*, in "Mol. Ecol.", 2013, vol. 22, n[o] 11, pp. 3179–90

[13] R. ASSAR, A. V. LEISEWITZ, A. GARCIA, N. C. INESTROSA, M. A. MONTECINO, D. J. SHERMAN. *Reusing and composing models of cell fate regulation of human bone precursor cells*, in "BioSystems",

April 2012, vol. 108, n⁰ 1-3, pp. 63-72 [*DOI :* 10.1016/J.BIOSYSTEMS.2012.01.008], https://hal.inria.fr/hal-00681022

[14] R. ASSAR, D. J. SHERMAN. *Implementing biological hybrid systems: Allowing composition and avoiding stiffness*, in "Applied Mathematics and Computation", August 2013, vol. 223, pp. 167–79, https://hal.inria.fr/hal-00853997

[15] R. ASSAR, F. VARGAS, D. J. SHERMAN. *Reconciling competing models: a case study of wine fermentation kinetics*, in "Algebraic and Numeric Biology 2010", Hagenberg, Austria, K. HORIMOTO, M. NAKATSUI, N. POPOV (editors), Springer, July 2010, vol. 6479, pp. 68–83 [*DOI :* 10.1007/978-3-642-28067-2_6], https://hal.inria.fr/inria-00541215

[16] M. BAKONYI, C. R. JOHNSON. *The Euclidean Distance Matrix Completion Problem*, in "SIAM J. Matrix Anal. App.", 1995, vol. 16, n⁰ 2, pp. 646-654

[17] E. J. CANDÈS, B. RECHT. *Exact Matrix Completion via Convex Optimization*, in "Found. Comput. Math.", 2009, vol. 9, pp. 717-772

[18] A. CARLSSON, J. YILMAZ, A. GREEN, S. STYMNE, P. HOFVANDER. *Replacing fossil oil with fresh oil - with what and for what?*, in "Eur J Lipid Sci Technol", 2011, vol. 113, n⁰ 7, pp. 812-831

[19] C. COMBES. *Parasitism: The Ecology and Evolution of Intimate Interactions*, University of Chicago Press, 2001

[20] P. GAYRAL, J. MELO-FERREIRA, S. GLEMIN. *Reference-free population genomics from next-generation transcriptome data and the vertebrate-invertebrate gap*, in "PLoS Genetic", 2013, vol. 9, n⁰ 4, e1003457

[21] M. HASIMOTO, S. HOSSAIN, A. AL MAMUN, K. MATSUZAKI, H. ARAI. *Docosahexaenoic acid: one molecule diverse functions*, in "Crit Rev Biotechnol.", Aug 2017, vol. 37, n⁰ 5, pp. 579-597, http://dx.doi.org/10.1080/07388551.2016.1207153

[22] L. LIBERTI, C. LAVOR, N. MACULAN, A. MUCHERINO. *Euclidean Distance Geometry and Applications*, in "SIAM review", 2014, vol. 56(1), pp. 3-69

[23] M. LYNCH. *Estimation of Nucleotide Diversity, Disequilibrium Coefficients, and Mutation Rates from High-Coverage Genome-Sequencing Projects*, in "Mol. Biol. Evol.", 2008, vol. 25, n⁰ 11, pp. 2409–19

[24] F. MORCILLO, D. CROS, N. BILLOTTE, G. NGANDO-EBONGUE, H. DOMONHÉDO, M. PIZOT, T. CUÉLLAR, S. ESPÉOUT, R. DHOUIB, F. BOURGIS, S. CLAVEROL, T. TRANBARGER, B. NOUY, V. ARONDEL. *Improving palm oil quality through identification and mapping of the lipase gene causing oil deterioration*, in "Nat Commun", 2013, vol. 4, 2160 p. , http://dx.doi.org/10.1038/ncomms3160

[25] R. E. RICKLEFS. *A comprehensive framework for global patterns in biodiversity*, in "Ecology Letters", 2004, vol. 7, n⁰ 1, pp. 1–15, http://dx.doi.org/10.1046/j.1461-0248.2003.00554.x

[26] F. RIMET, P. CHAUMEIL, F. KECK, L. KERMARREC, V. VASSELON, M. KAHLERT, A. FRANC, A. BOUCHEZ. *R-Syst::diatom: an open-access and curated barcode database for diatoms and freshwater*

*monitoring*, in "Database - The journal of Biological Databases and Curation", February 2016, vol. 2016 [*DOI :* 10.1093/DATABASE/BAW016], https://hal.inria.fr/hal-01426772

[27] S. T. ROWEIS, Z. GHAHRAMANI. *A unifying review of linear Gaussian Models*, in "Neural Computation", 1999, vol. 11, n⁰ 2, pp. 305–45

[28] L. K. SAUL, S. T. ROWEIS. *Think globally, fit locally: unsupervised learning of low dimensional manifolds*, in "Journal of Machine Learning Research", 2003, vol. 4, pp. 119–55

[29] D. W. THOMPSON. *On Growth and Form*, Cambridge University Press, 1917

[30] J. WANG. *Geometric structure of high-dimensional data and dimensionality reduction*, Springer & Higher Education Press, 2012