

DEVELOPMENT AND APPLICATION OF SINGLE-CELL ANALYSIS TOOLS FOR THE STUDY OF SYMPATRIC BACTERIAL POPULATIONS

Number of words: 27,663

Jasmine Heyse

Student number: 01205331

Promotors: Prof. dr. ir. Nico Boon
Prof. dr. Willem Waegeman

Tutors: ir. Benjamin Buysschaert
ir. Ruben Props
Peter Rubbens

Master's Dissertation submitted in fulfilment of the requirements for the degree of Master of Science in
Bioscience Engineering: Environmental Technology

Academic year: 2016 – 2017



De auteur en promotors geven de toelating deze scriptie voor consultatie beschikbaar te stellen en delen ervan te kopiëren voor persoonlijk gebruik. Elk ander gebruik valt onder de beperkingen van het auteursrecht, in het bijzonder met betrekking tot de verplichting uitdrukkelijk de bron te vermelden bij het aanhalen van resultaten uit deze scriptie.

The author and supervisors give the permission to use this thesis for consultation and to copy parts of it for personal use. Every other use is subject to the copyright laws, more specifically the source must be extensively specified when using from this thesis.

Ghent, June 2017

The promotors,

Prof. dr. ir. Nico Boon

Prof. dr. Willem Waegeman

The tutors,

ir. Benjamin Buysschaert

ir. Ruben Props

Peter Rubbens

The author,

Jasmine Heyse

Dankwoord

Met het afwerken van deze thesis kom ik aan het einde van vijf leuke en interessante maar ook uitdagende studiejaren. Het laatste jaar was daar allerm minst een uitzondering op. Het thesisjaar was zeker geen gemakkelijk jaar, maar desondanks vond ik het toch heel leuk en ben ik trots op het eindresultaat. Ik heb dit jaar heel veel bijgeleerd en bereikt. Dit alles zou nooit gelukt zijn zonder mijn drie tutors. Deze verdienen dan ook een oprecht bedankje. Ruben, ondanks de enkele duizenden kilometers tussenin kon ik altijd op jouw hulp rekenen en was geen vraag te veel voor jou. Zonder jouw aanstekelijk enthousiasme zou ik nooit zo verdiept geraakt zijn in de wereld van de microbiologie. Peter, jij stond altijd klaar om nieuwe input te geven waardoor ik het beste uit mijn data wou halen. Na een bemoedigend mailtje van jou had ik altijd opnieuw zin om er in te vliegen. Benjamin, jij hebt me met veel enthousiasme de wereld van het labo-onderzoek leren kennen, maar ook als het labowerk even niet meezat kon ik bij jou terecht. Bedankt voor jullie hulp dit jaar, ik heb heel graag met jullie samengewerkt.

Ik wil ook graag mijn twee promotoren, Prof. dr. ir. Nico Boon en Prof. dr. Willem Waegeman, bedanken. Bedankt voor jullie snelle, waardevolle en steeds positieve feedback. Ik ben blij dat ik in jullie beide onderzoeksdomeinen mocht delen.

Ook alle medewerkers van CMET verdienen oprecht een bedankje voor het creëren van een leuke sfeer, waardoor ik mij snel thuis voelde. Een speciaal bedankje gaat uit naar Tom, Cristina en Agathi. Tom, ik ben tientallen keren je bureau binnengevallen met een vraag over de flow cytometers. Bedankt dat jouw deur altijd open stond voor mij. Cristina and Agathi, thank you for your help and support during this last semester.

Graag wil ik Prof. Skirtach en Dmitry bedanken voor de hulp die mij geboden werd bij mijn experimenten met Raman spectroscopie.

Verder wil ik ook mijn medestudenten in CMET bedanken voor het gezelschap en de leuke babbeltjes tussenin, en dan in het bijzonder mijn twee partners in crime Lotte en Marjolein. Zonder jullie zou het nooit zo leuk geweest zijn tijdens de soms lange dagen in het flow cytometrie labo.

Tot slot wil ik graag nog mijn familie en vrienden bedanken. Mama, papa en Jasper bedankt voor jullie onvoorwaardelijke steun. Bedankt om me alle kansen te bieden en om me de vrijheid te geven om mijn eigen keuzes te maken. Bedankt, liefste checkers, voor de vijf leuke jaren, en al deze die nog zullen komen!

Contents

1	Literature study	1
1.1	Ecology of freshwater microbial communities	1
1.1.1	Microbial interactions	1
1.1.2	Freshwater ecosystems	3
1.2	Testing of ecological theories with synthetic ecosystems	5
1.2.1	Synthetic ecosystems	5
1.2.2	Bottlenecks	6
1.3	Phenotypes	7
1.4	Monitoring of microbial ecosystems	8
1.4.1	Bacterial growth and population density	8
1.4.2	Phenotypic community structure	10
1.4.2.1	Flow cytometry	10
1.4.2.2	Raman spectroscopy	15
1.4.2.3	Comparison	19
1.5	In silico communities	19
1.6	Objectives	21
2	Materials and methods	23
2.1	Bacterial strains	23
2.2	Molecular analysis	23
2.3	Microbial analysis	23
2.3.1	Bacterial growth	23
2.3.2	Flow cytometry	24
2.3.3	Raman spectroscopy	25

CONTENTS

2.4	Data analysis	25
2.4.1	Identification isolates	25
2.4.2	Bacterial growth	25
2.4.3	Flow cytometry	26
2.4.4	Raman spectroscopy	26
2.5	Experimental setups	27
2.5.1	Experiment 1	27
2.5.2	Experiment 2	28
2.5.3	Experiment 3	29
2.5.4	Experiment 4	30
3	Results	33
3.1	Experiment 1: Interactions between bacteria can lead to adjustment of their individual phenotypic diversities.	33
3.1.1	Phenotypic diversity assessment through flow cytometry	34
3.1.2	Prediction of relative abundances in the mixed culture	38
3.1.3	In silico approach	40
3.1.4	Phenotypic characterization through Raman spectroscopy	45
3.1.5	Difference in cellular composition	47
3.2	Experiment 2: Reversibility of the effect of interactions on the individual phenotypic diversities of the bacteria.	50
3.3	Experiment 3: Validation of the predicted relative abundances.	52
3.4	Experiment 4: Influence of carbon source diversity on phenotypic diversity.	53
4	Discussion	57
4.1	Hypothesis 1: Interactions between bacteria leads to an adjustment of their individual phenotypic diversities.	57
4.2	Hypothesis 2: Flow cytometry and Raman spectroscopy give complementary information regarding phenotypic community structure.	65

CONTENTS

4.3	Application of the in silico methodology to infer community composition	68
4.4	Experimental set-up	71
4.5	Conclusion and further perspectives	72
5	Appendix	1
5.1	Supplementary information data analysis	1
5.1.1	Bray-Curtis dissimilarity	1
5.1.2	Permutational multivariate analysis of variance using distance matrices (permanova)	2
5.1.3	PROcrustean Randomization TEST (PROTEST)	4
5.1.4	Random forest classifier	5
5.1.5	Randomized logistic regression (RLR)	6
5.1.6	Receiver operating characteristic (ROC)	8
5.1.7	Principal component analysis (PCA)	9
5.1.8	Principal coordinate analysis (PCoA)	10
5.1.9	t-Distributed Stochastic Neighbor Embedding (t-SNE)	10
5.2	Supplementary figures	14
5.2.1	Raman spectra nucleic acids	14
5.2.2	Lake isolates	15
5.3	Supplementary tables	16
5.3.1	Predicted relative abundances experiment 1	16
5.3.2	Criteria for cell classification	17
5.3.3	Relative abundances experiment 3	18

Table of Abbreviations

FCM	flow cytometry
FSC	forward scatter
LDA	linear discriminant analysis
LOOCV	leave-one-out cross-validation
PCA	principal component analysis
PCoA	principal coordinate analysis
PI	propidium iodide
RF	random forest
RLR	randomized logistic regression
SG	SYBR Green
SSC	side scatter
t-SNE	t-distributed stochastic neighbor embedding
OTU	operational taxonomic unit

Abstract

Ecosystems are often characterized by their diversity. The classic approach to evaluate diversity is to assess the genetic composition of a community. However, there is a broad level of diversity at an even finer scale, the phenotypic diversity. In literature, there is a growing interest in understanding how phenotypic diversity is manifesting itself and what its potential importance might be in both natural and engineered microbial ecosystems. This interest arises from the growing awareness that bacterial heterogeneity is an essential trait for many biological processes. Currently, our knowledge regarding factors that influence this phenotypic diversity is limited.

In this dissertation we have focused on the influence of microbial interactions on phenotypic diversity. Our results show that interactions between bacteria can lead to an adjustment of the individual phenotypic diversities of the interacting organisms. During this study we evaluated two techniques, flow cytometry, which had previously been used for phenotypic diversity estimation, and Raman spectroscopy, which had never been used for this purpose. Flow cytometry is a fast technique that gives information on cell-morphology and specific cell properties for which has been stained. Raman spectroscopy is a slower technique that gives a more holistic view on the molecular phenotypic traits. The potential of Raman spectroscopy as a tool to characterize phenotypic community structure was assessed. In addition, we evaluated its added value compared to the flow cytometric approach.

Furthermore, this study has been a part of research regarding development of the tools necessary for characterizing the community composition in synthetic ecology experiments. We have demonstrated that the experimental set-up and tools which were used during this study are suitable for studying both biotic and abiotic factors that might influence phenotypic community structure. This way we provided an experimental framework for further testing of ecological hypotheses regarding phenotypic diversity and microbial interactions.

Samenvatting

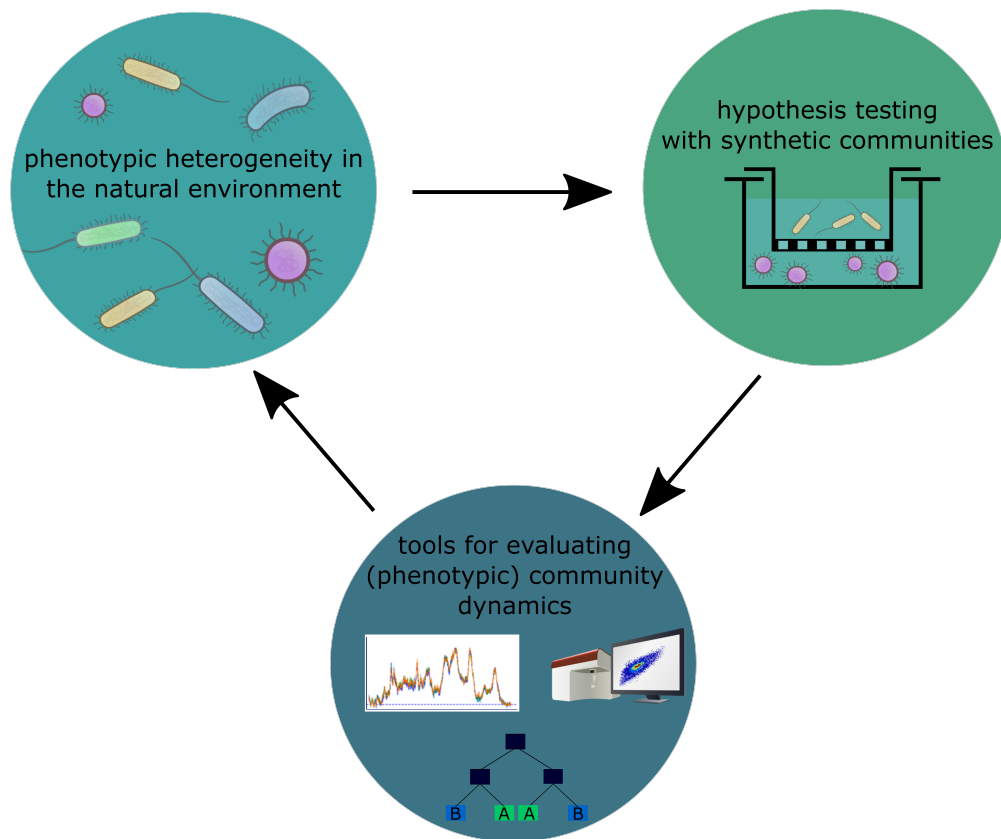
Ecosystemen worden doorgaans gekarakteriseerd aan de hand van hun diversiteit. Conventioneel wordt de diversiteit bepaald op basis van de taxonomische samenstelling van de gemeenschap. Naast taxonomische diversiteit is er echter nog een significante diversiteit binnen isogene populaties, namelijk de fenotypische diversiteit. In wetenschappelijke literatuur is er recentelijk een groeiende interesse in hoe deze fenotypische diversiteit tot stand komt en wat het belang van deze fenotypische diversiteit is in zowel natuurlijke als man-made ecosystemen. Deze interesse is ontstaan vanuit het groeiende besef dat bacteriële heterogeniteit van essentieel belang is in heel wat biologische processen. Desondanks is onze huidige kennis met betrekking tot de factoren die fenotypische diversiteit beïnvloeden is beperkt.

In deze masterthesis werd het effect van microbiële interacties op de fenotypische diversiteit van de interagerende bacteriën geëvalueerd. Onze resultaten tonen aan dat interacties tussen bacteriën kunnen leiden tot een aanpassing van hun individuele fenotypische diversiteiten. In deze studie werd gebruik gemaakt van twee technieken: flow cytometrie, een techniek die reeds gebruikt werd voor het inschatten van fenotypische diversiteit, en Raman spectroscopie, een techniek die nog niet gebruikt werd voor deze toepassing. Flow cytometrie is een snelle techniek die informatie geeft over cel-morfologie en specifieke eigenschappen waarvoor gestaind werd. Raman spectroscopie is een tragere techniek die een vollediger beeld van het moleculair fenotype kan geven. De mogelijkheid om Raman spectroscopie te gebruiken als een tool om fenotypische diversiteit in microbiële gemeenschappen te karakteriseren werd geëvalueerd. Bovendien werd de toegevoegde waarde in vergelijking met de flow cytometrie methodologie geëvalueerd.

Verder werd in deze studie ook nog actief meegewerkt aan de ontwikkeling van tools om de samenstelling van microbiële gemeenschappen in synthetische ecosystemen te achterhalen. Met dit onderzoek werd aangetoond dat de experimentele set-up en tools die gebruikt werden doorheen deze masterthesis geschikt zijn voor het bestuderen van zowel biotische als abiotische factoren die mogelijks de fenotypische diversiteit van microbiële gemeenschappen kunnen beïnvloeden. Op deze manier werd een experimenteel kader ontworpen dewelke kan gebruikt worden voor verder onderzoek naar fenotypische diversiteit en microbiële interacties.

Preface

Communities are often characterised and compared based on their diversity. When evaluating this diversity the classic approach is to evaluate the genetic composition of a community. However, there is a broad level of diversity at an even finer scale, the phenotypic diversity. A phenotype is a combination of all observable traits of an organism, and is assumed to be related to the functionality of the organism. Recently, there is an interest in understanding how this level of bacterial individuality is manifesting itself, which factors are influencing it and what its potential importance in both natural and artificial ecosystems might be.



To gain insights in factors that are causing and influencing this phenotypic heterogeneity, a first approach is to study it in a fixed environment, where potential influencing factors are being carefully controlled. These environments can be created through synthetic ecosystems. Sympatric bacterial populations are populations that are present in the same environment and therefore encounter each other on a regular basis, which implies they have the potential to interact. The goal of this study was to evaluate whether microbial interactions between these sympatric bacterial populations can lead to changes in phenotype, and at a broader scale, in phenotypic diversity of the community.

When we want to assess this fine scale diversity we need reliable tools to evaluate it. Since it is a phenomenon that is occurring at the level of individual bacteria, we applied two single-cell methods, flow cytometry and Raman spectroscopy. Both techniques give information regarding phenotypical traits, yet they have very different approaches and thus result in different information. When we want to filter out the relevant information from these analyses, we require a robust computational data-analysis.

1.1 Ecology of freshwater microbial communities

In nature, most organisms are not encountered as single cultures, but they are part of a larger association. A community is a naturally occurring group of organisms within a particular environment. These organisms and their interactions make up a functional unit, which is called the 'ecosystem'. The ecosystem consists of the biota (the members) and the habitat (the environment), and is often characterized by its diversity^[1].

The two components of diversity are richness and evenness. Richness indicates how many species are present, while evenness gives information on the relative abundances of the species. A low evenness indicates that a few species are dominating the system. Diversity influences the ecosystem functioning, but the exact mechanisms are unknown. Current theories propose two mechanisms. The first is that different species can use slightly different resources and therefore a more species-rich community results in a higher overall productivity of the ecosystem. Secondly, there is variation in the effect an individual species can have on ecosystem functioning. Some species will influence the community functionality more than others^[2]. Systems with higher evenness often have higher functionality and resistance to certain types of disturbances^[3]. It should be noted that diversity is not limited to species diversity alone, there exists a range of genetic and phenotypic diversity within single species^[1].

1.1.1 Microbial interactions

Microbial interactions occur both between members of different species as well as between members of the same species. The interactions result in the functioning and properties of the entire community. When studying organisms in mixed cultures, the observed community dynamics are often different from what would be expected when studying the different members in isolation. The

understanding of the rules and principles that govern these dynamics and community-properties is still limited. Therefore, understanding binary interactions is a necessary first step towards understanding the microbial community. However, observed dynamics in studies with more than two genotypes can still differ from what would be expected based on known binary interactions^[4]. Interactions might involve physical contact, such as gene exchange, or there might be a uni- or bidirectional exchange of small chemicals, such as in case of metabolic interdependencies. Organisms can also provide a function for other organisms, such as the degradation of antibiotics^[5].

Interactions can be split up into two categories, the active and the passive interactions. One organism might be producing some waste product, which can potentially be used as a food source for a second organism. This interaction is passive since the first organism is not intentionally promoting the growth of the second organism. In case an organism is intentionally investing resources in his metabolic process or behavior in order to influence another organism, the interaction is active^[4].

Classification of the interactions can also be based on the resulting effect for the interacting organisms. The interaction might be beneficial, neutral or detrimental for the involved organisms. Depending on this, the interactions are indicated with different terms. 'Commensalistic' interactions are interactions where one organism is benefiting from the relationship, while the other neither has an advantage nor a disadvantage. An example of a commensalistic interaction can be cross-feeding of one organisms by the waste products of another organism, or depletion of an antibiotic by an antibiotic-degrading organism, which then allows non-resistant organisms to survive. 'Mutualism' describes relationships where two organisms are both benefiting. For example, two organisms might both actively invest energy in their metabolisms in order to produce products that can be used by their partner^[4]. Mutualistic interactions can be disturbed by so-called 'cheaters'. These are individuals that make use of the benefits of the interactions in the community, but do not contribute by producing resources themselves. This might lead to a breakdown of the mutualistic interaction^[6]. The cooperating organisms can protect their interaction with different types of anti-cheating mechanisms: targeted benefit and targeted punishment. In targeted benefit the cooperating organisms get access to benefits, which might vary according to their contributions to the interaction. In the targeted punishment the cooperative organisms will actively punish the cheaters, for example by producing a toxin that does not affect the cooperative organisms themselves^[7].

A negative kind of interaction is competition. Exploitative competition, which is competition for resources and space, can potentially have an important influence on the shape of the community, mostly by causing selective extinction. This is referred to as the 'competitive exclusion principle', which states that two species which are competing for the same limiting resource can not exist

together. In other words, there is a limit on how many species can exist in the community before a certain niche becomes saturated. In some cases species evolve towards different niches in order to coexist, this is called 'divergent co-evolution'. Niche differentiation can occur in space and time, but most often manifests itself as a morphological differentiation^[8]. Competition can be a barrier for new organisms to enter the system^[9].

Not only the community members decide the interactions, the environment also has an influence. For example, in a homogeneous environment it is more likely that competition will lead to extinction of certain members in comparison to a more heterogeneous environment^[8]. Sometimes there can be a stable coexistence of similar species because of spatial organization of connections and interactions^[6]. The chemical composition of the environment also has an influence on the interactions. Interactions can be induced or changed, not by changing the interacting microorganisms themselves, but by changing their growing medium^[5].

1.1.2 Freshwater ecosystems

Bacterial communities have relatively low cell densities in natural freshwater environments, typically 10^5 - 10^6 cells/mL^[10]. Nevertheless, they are important members of lake ecosystems, where they play a key role in the transformation and cycling of biologically active elements. Still, the bacterial taxa that play the most prominent role in these ecosystems remain relatively unknown^[11]. The current knowledge on aquatic microbial diversity is mainly based on the analysis of 16S rRNA gene sequences that were produced by high-throughput sequencing of environmental DNA or on the knowledge derived from databases that were established using rich solid media^[12]. Despite that the bacterial community present in lakes can vary enormously between different lakes, there are some groups of freshwater bacteria who are widely distributed^[13]. Betaproteobacteria are the most abundant bacteria in the upper water column of lakes, especially the two genera *Polynucleobacter* and *Limnohabitans*. The class of Betaproteobacteria is split up into seven lineages: betI to betVII. The betI lineage is further split up in two clades betIA and betIB, which is equivalent to splitting it up in the genera *Limnohabitans* and *Rhodoferrax*. *Limnohabitans* species are ubiquitous in freshwater ecosystems and typically occur with high relative abundances, which makes it an environmentally important group^[11]. The genus *Limnohabitans* was only recently established, in 2010^[14]. Phylogenetic analysis has revealed the presence of different tribes within the genus, which are indicated as Lhab-A1, Lhab-A2, Lhab-A3 and Lhab-A4^[11]. However, some studies have suggested other subdivisions of the genus^[15].

Within the genus of *Limnohabitans*, there is a high level of diversity, both morphologically and physiologically. The large variability in cell sizes and morphologies is illustrated in Figure 1.1. Cell sizes can be linked to differences in grazing pressures^[16]. One of the factors that influence the composition of the bacterioplankton in lakes is the trophic status, since this largely determines the available niches and the composition of zooplankton that graze on the bacterioplankton^[13]. Another important factor in shaping the community composition of bacteria in freshwater lakes is pH. It is unknown whether pH directly influences the bacterial community or whether it indirectly influences other growthfactors^[17]. *Limnohabitans* species have shown a large diversity in pH-preference, ranging from slightly acidic to alkaline waters, but with a more frequent preference for alkaline habitats^[18;19]. Since morphological and physiological differences indicate differences in ecology, this observed diversity suggests that *Limnohabitans* might inhabit a wide variety of niches in standing freshwater ecosystems^[15], explaining why species of this genus are found so frequently in different types of freshwater ecosystems^[18]. The clear habitat preferences that have been found for different *Limnohabitans* species indicate that their omnipresence cannot be explained by the species to be ‘generalists’, but rather indicates a high level of specialization for different species^[18].

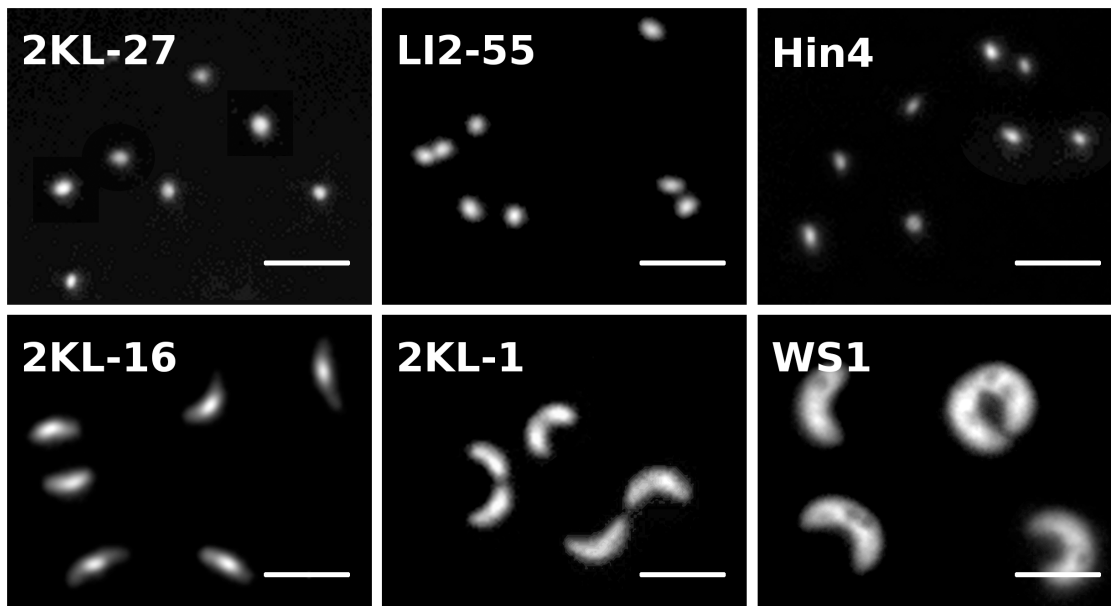


Figure 1.1: Different sizes and morphologies of *Limnohabitans* isolates. The codes refer to identification labels that were assigned to the isolates during the study of Kasalický *et al.*. The scale bar represents 2 μm ^[15].

1.2 Testing of ecological theories with synthetic ecosystems

1.2.1 Synthetic ecosystems

Artificial communities, commonly referred to as ‘synthetic ecosystems’, are used to study microbial ecosystems. Communities with the same key properties as a natural ecosystem of interest can be created. A synthetic ecosystem consists of a selected set of species under specified conditions. They have a reduced complexity in comparison to natural communities, which offers simplicity. Next to their simplicity, also their controllability is an advantage^[20]. These artificial ecosystems can be used as a model to study particular processes and their influencing factors. Information about structure, evolution and functioning of the microbial communities can be obtained^[21]. However, these artificial ecosystems should not be seen as miniature versions of the actual ecosystem, but rather as a way to test ecological theories to better understand the rules of nature and to answer questions that are difficult to study directly in the field^[22]. Researchers should be aware of the fact that lab-observations do not always match the dynamics in natural systems exactly^[23]. Synthetic ecosystems allow to study organisms in mixed microbial communities rather than in pure cultures, which might give unexpected insights since the behavior of an organism in a pure culture can be different from its behavior in a mixed microbial community^[20].

A specific set-up for synthetic ecosystems are co-cultures. The principle of such a system is that two or more populations are grown together with some degree of contact between them^[24]. These systems can be used to cultivate organisms that are not easily monocultured because they need the presence of certain other organisms^[25]. They are also often used as tissue models for medical applications such as drug testing or host-microbiome interaction studies^[26]. Co-culture systems appear in many different forms including gel encapsulated cells, Petri dishes and membrane separated systems such as transwells (Figure 1.2)^[24].

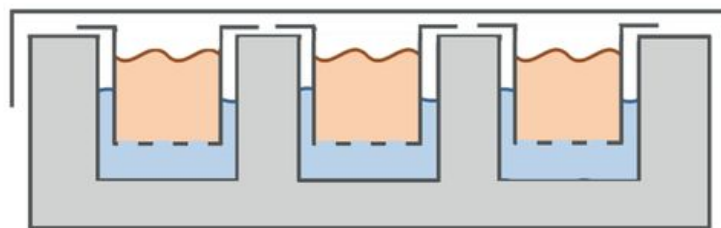


Figure 1.2: Transwell system for co-culturing. The basale and apicale phase are separated using a semi-permeable membrane^[24].

Research based on synthetic ecosystems is sometimes criticized for being too simplistic because the degree of complexity is much smaller in comparison to that of natural systems. This complexity is determined by the researchers instead of natural processes. The artificial nature of these ecosystems is also criticized. Researchers have control over the inputs and constituents of the system, such as environmental conditions and initial abundance of the organisms, which have been shown to be important drivers for ecosystem structure and functioning. However, this low complexity allows to exclude certain influencing factors or to assess only influencing factors of interest, which is a necessity for testing ecological theories^[22]. Once the experiment has started, the dynamics that are observed are not created or directly controlled by the experimenter. And the species which are combined to form these artificial ecosystems are mostly species that are found together in nature^[22]. Experiments in the full natural context would often be too complex to allow clear conclusions^[27]. Micro-organisms have ideal properties for synthetic ecology studies as they have short generation times and small sizes. The short generation times allow to study processes that are typically occurring over multiple generations, such as co-evolution, while the small sizes make it practically possible to work with replicates^[22].

1.2.2 Bottlenecks

When synthetic ecosystems are used to study fundamental principles of nature, there are some difficulties. The change in community composition over time is often of interest. For example, one wants to know which species will dominate the community after a certain period of time or under certain conditions. A tool to evaluate this change in community composition is needed. Classically, community composition is evaluated based on plate counting methods or molecular techniques. However, these techniques have some drawbacks. They can be time consuming, which is not feasible when one wants to follow up a change in community on relatively small time intervals. Additionally, significant bias can be introduced in case cultivation-based methods are used^[25]. Another difficulty is that molecular techniques are invasive, which implies that often only one sample can be taken as most experimental set-ups are relatively small.

Another problem encountered when working with synthetic ecosystems is that the growth kinetics of the individuals in the community is known for pure cultures, while it will most likely be different in a mixed community. But this growth is mostly not known, nor is it easily determined^[28].

1.3 Phenotypes

The phenotype of an organism refers to its observable traits. A phenotype is mostly defined as the result of interaction between the genotype and the environment^[29]. Over the years the concept of phenotype has been expanded from morphological traits, to also ‘molecular phenotypic traits’, such as protein contents and the mRNA levels of genes. The current definition of phenotype holds all ‘observable’ characteristics^[30]. Different types of assays have become available to evaluate these traits. These are called ‘phenotypic arrays’. Mostly they are focused on evaluating respiratory or growth characteristics of a certain species or community under different conditions (i.e. different growing media, temperatures, etc.)^[31].

There is a growing interest in the influence of phenotypic traits and phenotypic heterogeneity on interactions, communities and ecosystem functions. In light of this, ‘phenotypic plasticity’ is often of interest. Phenotypic plasticity is the ability of an organism to change its phenotype as an adaptive response to a changing environment. This can relate to both changes in behaviour and in morphology. A large plasticity indicates there can be large changes in phenotype in response to the environment. Phenotypic plasticity has been reported to influence competition between species and the possibility of coexistence of different species related to this competition^[32].

Phenotypic heterogeneity occurs in sympatric isogenic populations, i.e. isogenic populations which coexist in a given habitat^[33]. This heterogeneity is explained by a few processes. One of those is that gene expression and other cellular processes are low rate processes regulated by small molecules, which are present in small amounts in cells. These molecules are unevenly distributed within the cell, and during the cell division they will be distributed unequally in both daughter cells. This unequal distribution results in stochastic variation in cellular reactions and gene expression, which is called ‘biological noise’^[34]. Another factor that can cause heterogeneity at this level are periodical changes in the cellular function, such as cell cycles or switches between different metabolic processes^[34]. Other factors, such as cell to cell interactions are also in play. Cells make ‘decisions’ regarding their gene expression, but it is very likely that these decisions are influenced by other members of the population^[34]. This diversity within isogenic populations can have implications on both functionality and survival. There can be a division of labour between phenotypes, which increases overall productivity of the species^[34]. Some phenotypes might be able to deal better with a certain change in conditions. In this way the heterogeneity might allow subpopulations of the species to persist during changing conditions, and thus increase the survival probability of the species^[35]. The discovery of this heterogeneity, which is by definition independent of genetics

or environmental conditions, has led to a renewed view on microbial individuality^[36].

Some new questions will need to be answered, such as how important this type of diversity is in natural environments and whether this diversity is negligible in comparison to genetic diversity and diversity caused by environmental gradients^[34]. This level of diversity, has some of the same key properties as genetic diversity, such as that it can offer resilience and functionality. However there are also some fundamental differences. Phenotypes can not be depleted by continuously changing environmental conditions, while genotypes can^[34]. The disappearing phenotypes can be replenished from other phenotypes of that genotype^[34]. Division of labour between phenotypes where one phenotype is providing something for others but does not receive anything in return is possible, whereas this would not be possible for different genotypes^[34].

1.4 Monitoring of microbial ecosystems

1.4.1 Bacterial growth and population density

Bacteria reproduce by binary fission. This means that each cell grows larger and then divides into two smaller daughter cells, which then in their turn grow and divide. Provided that no mutations occur during the cell division, the daughter cells are genetically identical to the original cell and form isogenic populations^[37]. The time it takes for a bacterial cell to divide, is called the generation time. So each generation time, the amount of bacterial cells is doubled. These generation times can differ from a few minutes to several hours^[38]. For optimal growth certain nutritional and physical factors, such as pH and temperature, must be met. Nutritional factors can be split up in macro- and micro-nutrients and include amongst others nitrogen, phosphorus, carbon, sulphur and water^[39]. When growing bacteria in laboratory conditions, these nutritional factors are provided via the growth medium. Certain environmental conditions can be simulated by choosing a growth medium with specific limiting factors^[40].

The Monod equation gives the empirical relationship between the specific growth rate (μ) and the concentration of the substrate that is limiting for growth (S) (Equation 1.1). The specific growth rate is the relative amount of cells that is produced per amount of cells that is present per time. The equation has two empirical parameters: the maximum specific growth rate (μ_{\max}) and the affinity constant (K_s). This affinity constant is the substrate concentration S when μ is half of μ_{\max} ^[41].

$$\mu = \mu_{\max} \frac{S}{K_s + S} \quad (1.1)$$

Bacterial growth is classically presented in a graph that shows time on the x-axis and the logarithm of the number of cells on the y-axis (Figure 1.3). This curve can be divided into four phases: the lag phase, the log or exponential phase, the stationary phase and the death or decline phase. During the lag phase there is none or very slow growth because the organisms are adapting to their culture conditions. In the log or exponential phase, the cells double at optimal growth rates. In the stationary phase, cell numbers stabilize as the medium gets depleted and metabolites start to accumulate. There is as much growth as there is death during this stationary phase. Finally, during the death phase there is a net loss of viable cells^[42]. Of course, this is a theoretical model, and in practice several exceptions to this classic growth curve have been observed, such as growth curves with two exponential phases, which are called ‘diauxic growth curves’^[43].

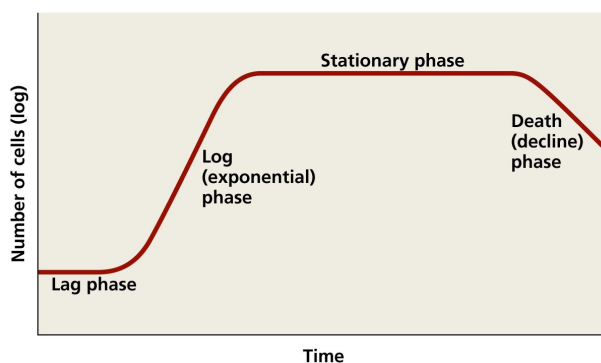


Figure 1.3: Typical shape of the bacterial growth curve. Four phases can be distinguished: lag phase, log phase, stationary phase and death phase. A logarithmic scale is used on the y-axis since bacterial growth is exponential^[44].

The properties of cells are not constant during growth, size and shape of the cells change during the different growth phases^[45]. In other words, bacteria exhibit different phenotypes throughout the population growth. This implies that mass increase and increase in cell density are not directly related. When growing conditions become unfavorable, a large surface-to-volume-ratio is more feasible. And thus during the stationary phase, bacteria will be smaller and less spherical, while the opposite will be true during the exponential phase^[46]. Cell size increases during the lag phase, reaches a maximum during the log phase, and eventually decreases again in the stationary phase^[47]. Also chemical composition of the cell, for example DNA content, is changing during the growth process. Sometimes size and protein content of cells can remain constant during exponential growth, however they can also change because the medium in which they are growing is changing over time, due to the use of nutritional compounds and the accumulation of metabolites^[29].

1.4.2 Phenotypic community structure

Phenotypic diversity has only been partly explored in natural systems. One of the reasons for this is that it is difficult to assess single species diversity within mixed communities^[1]. Two laser-based methods that are suitable for assessing phenotypes are flow cytometry and Raman spectroscopy.

1.4.2.1 Flow cytometry

Flow cytometry (FCM) is a laser-based technology that analyses single particles, usually cells, by sending them through a beam of light by a fluid stream. The general set-up is given in Figure 1.4. In order to be able to interrogate every particle separately, there is a fluidics system that transports the particles through the laser beam one by one. This fluidics system consists of a stream of sheath fluid that directs the sample stream to the center of the laser beam. This process is known as ‘hydrodynamic focusing’^[48].

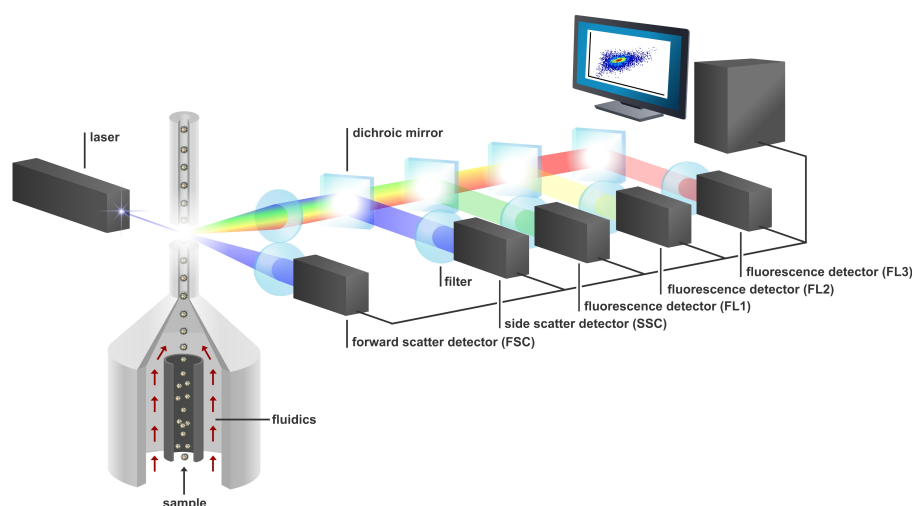


Figure 1.4: General set-up of a flow cytometer. The sample stream is directed towards the laser by a stream of sheath fluid. The particles of the sample are sent through the laser beam one by one. Scattered light and fluorescence are collected by a system of dichroic mirrors and photomultipliers^[49].

Two types of light can be detected in the flow cytometer, that is scattered light and fluorescence. While the particles pass through the beam, they will scatter light, which can be measured by a detector in front of the light as forward scatter (FSC) and by one or several detectors to the side as side scatter (SSC). Particles can also be stained with fluorescent stains prior to analysis. These stains will be excited by the laser, which will cause them to emit fluorescent light. The fluorescence intensity is proportional to the amount of stain that is bound to the particle, thus

there is a stoichiometric staining^[50]. The scattered and fluorescent light are collected by a system of dichroic mirrors and photomultipliers^[51]. In general, the scattered light provides information about the basic characteristics of the cells (e.g. size, shape and surface properties). The fluorescent data provides additional features which can be used to characterize the bacteria, distinguish the bacteria from abiotic particles and indicate cell viability and vitality^[52].

Depending on the fluorochrome that is used for staining, information about different cell characteristics such as membrane integrity, intracellular pH and metabolic activity can be obtained^[53]. The multitude of available stains gives FCM its broad applications in the fields of industrial biotechnology, immunology, pharmaceuticals, microbial ecological research, etc.. Flow cytometers that are equipped with multiple lasers and detectors allow for assessing multiple parameters simultaneously^[51]. However, when multiple fluorescent markers are used, optical crosstalk poses a significant issue. Cross-talk or spill-over occurs when fluorochromes have overlapping emission spectra, and thus the detector that is intended for one fluorochrome will also measure signals coming from the other fluorochrome^[54]. In some cases there is no need for staining because of autofluorescent properties of certain components. Chlorophyll a, a pigment present in for example cyanobacteria, emits a red fluorescence, and thus cyanobacteria can be detected without any staining^[55]. Also bacteria with fluorescent labels such as a *gfp*-label (green fluorescent protein) can be used as such.

An often used application of flow cytometry is live/dead staining. For this, mostly a combination of SYBR green I (SG) and propidium iodide (PI) is used, however other options exist. Both SG and PI bind to nucleic acids and thus stain DNA and RNA. SG can permeate into all cells, which makes it a useful stain for total cell counts, while PI can only enter cells with a damaged cytoplasmic membrane. This way intact cells will only be stained by SG, while damaged cells will be stained by both SG and PI. The emission spectrum of SG overlaps with the absorption spectrum of PI. Therefore, the emission of SG is used as excitation energy for PI, which then in turn emits red fluorescence while the SG fluorescence is quenched^[56]. This difference in fluorescence signals results in the fact that two groups can be distinguished by the flow cytometer, one group with intact cells and one with damaged cells. Membrane integrity is widely accepted as a criterion to evaluate cell viability^[57].

Speed is a big advantage offered by FCM^[51], up to 50,000 cells per second can be analysed^[58]. The ability to use automated sample loading allows fast analysis of large amounts of samples. When using autosamplers, care should be taken that the delay until the measurement of the last sample is not too large, since this might cause bleaching of the stain. Sedimentation of larger cells could also occur, which would have large effects on the results of the measurement^[59].

In contrast to conventional heterotrophic plate counting methods, microbial cells can be detected by FCM irrespective of their cultivability. Since it is estimated that less than 1% of bacteria in aquatic environments can be cultivated by traditional plate counting techniques^[60], this is an important benefit. Another benefit is that even in systems with low cell densities large amounts of cells can be analysed^[51]. This is particularly beneficial when one wants to study samples of natural freshwater systems, since cell densities are typically low there.

Data generated by FCM analysis is generally presented in single parameter histograms or in two-dimensional scatterplots. These two-dimensional plots are often referred to as ‘fingerprints’ (Figure 1.5). Abiotic particles, such as crystals or dust, can scatter light and potentially bind with the stains or possess autofluorescent properties. Hence they might interfere with the signals that truly originate from the cells^[61]. When one wants to interpret the data correctly, these datapoints need to be removed. This is done by so-called ‘gating’, illustrated in Figure 1.5. Gating is selecting the area of interest and is mostly performed based on signal intensities of two fluorescent parameters^[51]. Multiple gates can be created to differentiate between subpopulations of bacteria in a sample. Gating is mostly performed manually by the operator and is therefore subjective, which implies that care should be taken when comparing data of different datasets. Gates can also be created using clustering algorithms^[62]. Next to gating, there is another parameter which makes it difficult to compare between separate datasets: the operator can change the voltages over the photomultipliers in order to tune the detectors. Higher voltages will cause more electrons to be generated per event, which implies also smaller events will be detected well by the device.

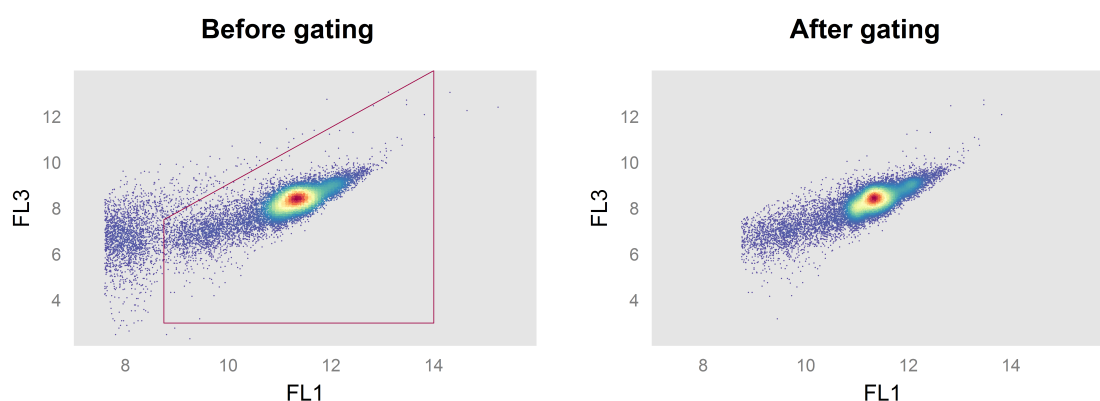


Figure 1.5: Illustration of the gating step. **Left:** Data before gating. The gate is manually created based on signal intensities of two fluorescent parameters, here denoted as FL1 and FL3, in order to isolate the cellular information from the background. **Right:** Isolated cellular information after gating. The colour intensity is proportional to the log-scaled density of the events.

As stated above, FCM can be used to gain information on phenotypic traits (i.e. size, shape, nucleic acid content etc.) of single cells. Recently, a new method has been developed to assess phenotypic community structure of bacterial populations based on flow cytometric fingerprinting^[63]. The pipeline will be explained here.

A first step when processing the FCM data is applying a transformation, such as a hyperbolic arcsin function. The hyperbolic arcsin will transform low values approximately linearly and high values approximately logarithmically. The reason for this transformation is that there is a broad spread in values of the flow cytometric parameters (i.e. scatter and fluorescence). By applying a transformation the variance in the data will be reduced, which will make the density estimation more accurate. The next step is denoising the data in order to extract only the signals that are resulting from bacteria, and removing all signals originating from the background (i.e. particles, salts, etc.). For this, a gating step is performed (Figure 1.5). Next, the data is normalized to the [0,1] interval by dividing each parameter by the maximum FL1 (i.e. the first fluorescence channel) intensity value over the data set. The reason for this is to make a bandwidth of 0.01 appropriate during the density estimation. The next step is to calculate the phenotypic community structure. First, a discretization step is performed by applying a 128x128 binning grid (Figure 1.6 A and B). This grid is applied for each of the binary parameter combinations. For example, when using two scatter and two fluorescence parameters for characterizing the phenotypic community structure, this would result in 6 binary parameter combinations. During the discretization step the operational phenotypes (bins) are being defined. In each of the bins a Gaussian distribution is fitted by Kernel density estimation, using a bandwidth of 0.01. Finally, all density estimations are summed, leading to the density estimation of the community (Figure 1.6 C). The density values for all bins are then concatenated into a 1D-vector, which is called the 'phenotypic fingerprint'(Figure 1.6 D). This phenotypic fingerprint has a structure similar to the output of sequencing pipelines, containing the probability that a cell is present in each bin. And thus established ecological indices can be calculated from it.

To evaluate alpha-diversity (i.e. within sample diversity) the first three Hill numbers are calculated: D_0 (species richness), D_1 (the exponential of Shannon entropy) and D_2 (the inverse of Simpsons index). Formulas are given in Table 1.1. The larger the Hill order, the lesser rare species are taken into account. Hill numbers can be interpreted in terms of 'effective number of species', which represent the number of equally abundant species required to generate an identical diversity as the one of the microbial community under study. Thus all Hill numbers have the same units as species richness^[64]. They also obey the doubling property which states that when N equally diverse groups

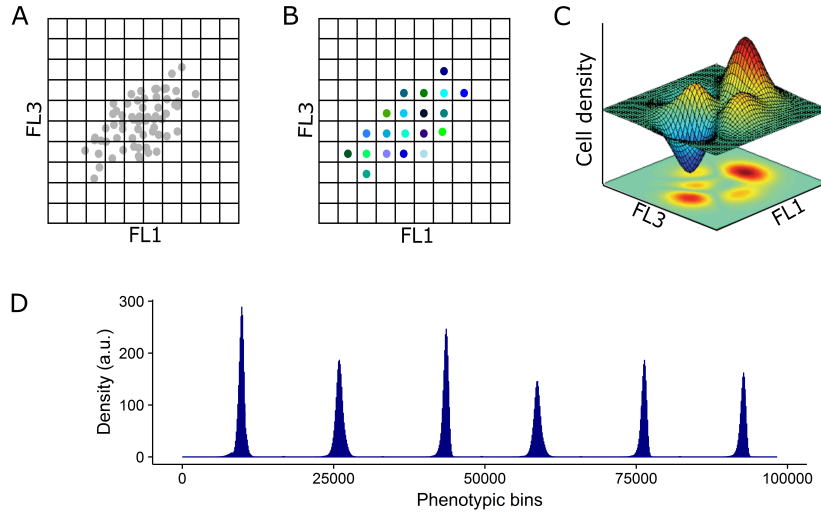


Figure 1.6: Conceptual illustration of the calculation of the phenotypic fingerprint. **A:** The data is transformed, denoised and normalized. **B:** A binning grid of 128x128 is applied on each of the binary parameter combinations. This binning grid defines the phenotypes. **C:** For each of the bins a distribution is fitted by Kernel density estimation. All these density estimations are summed to create the phenotypic community structure. **D:** The values for each of the bins are then concatenated into a 1D-vector.

which have no species in common are merged, the merged community has a diversity that is equal to N times the diversity of one of the groups^[64]. This makes Hill numbers easy to interpret.

Table 1.1: Formulas of the order-based Hill numbers, where p_i is the relative abundance of bin i and S is the number of non-empty bins in the phenotypic fingerprint.

Hill order (q)	Diversity metric (D_q)
0	$D_0 = S$
1	$D_1 = e^{-\sum_{i=1}^S p_i \ln(p_i)}$
2	$D_2 = \frac{1}{\sum_{i=1}^S p_i^2}$

Beta-diversity (i.e. between sample diversity) is evaluated using the Bray-Curtis dissimilarity. The Bray-Curtis dissimilarity for two samples A and B is given in Equation 1.2, where p_{Ai} is the abundance of bin i in sample A and S is the total number of non-empty bins in the phenotypic fingerprints of both A and B. This metric takes a value of zero in case of two identical samples, and a value of one in case of two samples which have no bins in common^[65] (for more information

on Bray-Curtis dissimilarity see Appendix 5.1.1).

$$D_{AB} = \frac{\sum_{i=1}^S |p_{Ai} - p_{Bi}|}{\sum_{i=1}^S (p_{Ai} + p_{Bi})}. \quad (1.2)$$

1.4.2.2 Raman spectroscopy

Raman spectroscopy is a laser-based technology that assesses the chemical composition of a sample. Molecules consist of atoms that are elastically bound to each other. These elastic bonds have periodical motions, the vibrational modes, which include bending and stretching. Raman spectroscopy is a form of molecular spectroscopy that studies these vibrational modes. It is based on scattering of monochromatic light. When monochromatic light is sent to a sample, this light will be partly scattered. Scattering can be split up in two types: elastic and inelastic scattering (Figure 1.7). The majority of the scattered light has the same wavelength as the incoming light, this scatter is referred to as elastic or Rayleigh scatter. When there is a change in wavelength, this scatter is referred to as inelastic scatter or Raman scatter. There are two types of Raman scattering: Stokes and anti-Stokes. In case the photons get a lower energy level, and thus a higher wavelength, compared to the incoming photons, the scatter is called Stokes scatter. While an increase in energy level of the photon as compared to the incoming photons is called anti-Stokes scattering. Molecules that are initially in the ground state give rise to Stokes scattering, while molecules that are initially in an excited vibrational state give rise to anti-Stokes scattering. At ambient temperatures, more molecules are in their ground states, and therefore Stokes scattering is more intense than anti-Stokes scattering^[66;67].

The shift in energy is caused by interaction of the monochromatic laser light with the vibrational states of the molecules that are present in the sample^[69]. The vibrational state of the molecules causes an electric field. The polarizability is the ease with which the electron cloud can be changed in shape, size or orientation as a response to an external electric field. A photon will also generate an electromagnetic field, with a frequency which is proportional to its energy level. This electromagnetic field, generated by the photon's energy state, polarizes and vibrates the electromagnetic field of the molecule. This is called an induced dipole. The photon can lose some of its energy by causing this induced dipole or gain some energy in case of a constructive interaction with the induced dipole. So this shift in energy, and thus wavelength, of the photon is the result of field-interactions with the vibrational energy of the molecule^[66]. The vibrational modes of a molecule are dependent on mass of the atoms and their geometric arrangement, the nature of the chemical

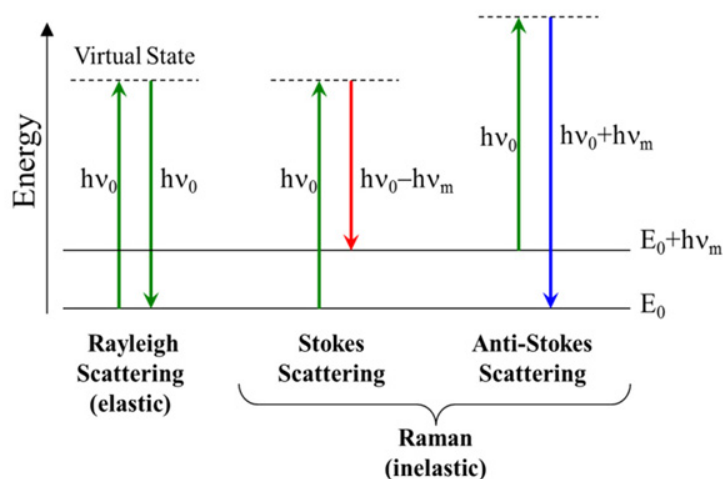


Figure 1.7: Illustration of the different types of scattering. Horizontal lines represent different vibrational states of the molecular bond. Incoming photons have an energy equal to $h\nu_0$, with ν_0 the incident photon frequency and h is Planck's constant. In case of elastic or Rayleigh scattering, the scattered photons have the same energy as the incoming photons. In case of Stokes scattering, the scattered photons have a lower energy and the molecule gets a higher energy, while the opposite is true for anti-Stokes scattering^[68].

bonds and their motions. Raman scattering can thus be used as information source for structure and properties of the molecules present in a sample^[67]. Next to that, the Raman signal is also dependent on concentration of the molecules, which makes quantification of chemical compounds possible^[66].

As mentioned above, there needs to be a change in polarizability to create an induced dipole and thus in order to be able to get Raman scattering. The larger the polarizability, the larger the Raman-effect. This implies that only certain vibrational modes will be Raman-active, and thus not all compounds can be detected using Raman spectroscopy. For example, water has a very low Raman activity, which results in the fact that the presence of water will not interfere with the Raman signal^[70].

Raman scattering is inherently a very weak process, only 1 in $10^6 - 10^8$ incident photons are Raman scattered^[71]. The Raman spectrum contains the intensity of the scatter signal as a function of the energy difference with the incident light, expressed in terms of the so-called 'wavenumber'. This wavenumber is calculated based on the wavelength of the incoming photons (λ_0) and of the scattered photons (λ_1), the relation is given in Equation 1.3. Wavenumbers are typically expressed in cm^{-1} . An important remark here is that the wavenumber is a relative value and thus depending on the wavelength of the laser. This implies that spectra obtained with different lasers cannot

be compared directly^[72]. The Raman-peaks are indicative for different chemical bonds and their vibrations^[71].

$$\Delta w = \frac{1}{\lambda_0} - \frac{1}{\lambda_1} \quad (1.3)$$

The general setup of a Raman spectroscope is given in Figure 1.8. A laser is used as monochromatic light source. This laser beam is sent to the sample through a monochromator which ensures the incident light has only one particular wavelength. The incident light can interact with the sample and the resulting scattered light is focused through a lens. The scattered light is usually observed in the direction perpendicular to the laser beam. The signal is sent through band pass filters in order to remove all the photons that have resulted from elastic scattering. The resulting light consists of inelastically scattered photons and is sent to a grating. This grating will reflect each of the scattered wavelengths under a slightly different angle. In this way each wavelength is collected on a separate spot on the detector. The photons are typically measured at lower energy than the input light, which means the Stokes scattering is detected. However, also the anti-Stokes could be used^[69].

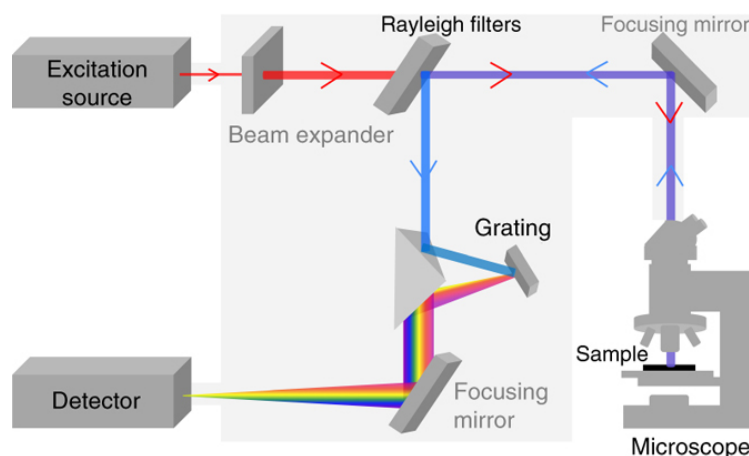


Figure 1.8: The laser beam is sent through a monochromator and focussed on the sample. The laser light interacts with the sample and the resulting scatter is collected through a lens. The light is sent through a set of filters to remove Rayleigh scatter. After grating, each wavelength is detected by the detector^[73].

Raman spectroscopy has applications in a variety of fields, but is mainly used for material identification, for example to identify mineral composition in geology, or to assess purity of polymers and pharmaceuticals^[74]. It is also frequently used for applications in life sciences such as single-cell studies. The Raman spectrum of a single cell is a combination of the spectra of all the different compounds that make up this cell (e.g. proteins, nucleic acids, fatty acids, etc.). This results in

a very complex spectrum, which can be interpreted as a chemical fingerprint of the cell^[75] (Figure 1.9). The spectrum can be used for phenotypic characterization of the cells^[76]. The spectra are used as an identification tool to distinguish between bacterial species or for the interpretation of the presence of biomolecules. For the latter, a database of the spectra of the biomolecules is needed. Some studies have been focussing on creating databases with reference spectra^[77]. Because of the complex nature of the Raman cell spectrum, interpreting it is not straightforward^[77]. In case one wants to use the spectra as an identification tool, one should take into account that the spectrum is not only dependent on the species or strain that is evaluated, but also on the metabolic history of the organism^[71].

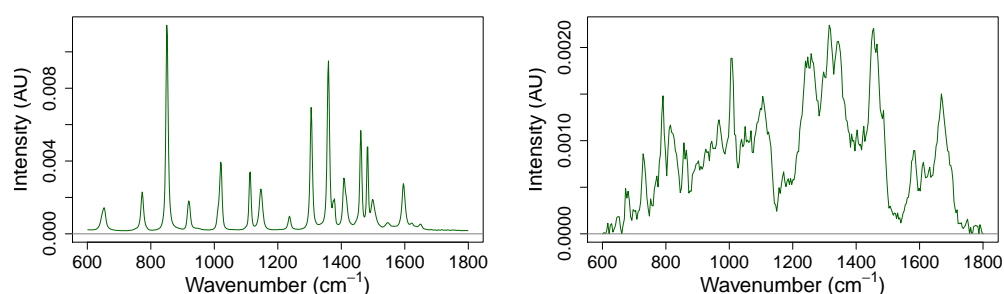


Figure 1.9: Left: Raman spectrum of alanine, an amino acid. **Right:** Raman spectrum of single cell.

Even simple natural microbial communities consist of a large amount of species, and it is not straightforward to determine which species are playing an active role in the community or are performing key functions^[21]. One goal in environmental microbiology is to link particular functionalities or activities to particular organisms. A way to do this could be to use Raman spectroscopy in combination with stable isotope tracers^[71]. This method has successfully been applied by tracking of heavy water incorporation to evaluate cellular activity at the single cell level^[78].

Like flow cytometry, Raman spectroscopy is cultivation independent. Thus also viable but non culturable cells can be analyzed with Raman spectroscopy, which is an interesting property since most of bacterial species occurring in nature cannot be grown easily under laboratory conditions^[79]. Another advantage of the technique is that it does not require any stains or specific reagents^[75].

The differences between bacterial spectra can be very subtle and therefore robust computational data-analysis is highly required in order to filter out the relevant information. The spectra require some pre-processing. There is a wide variety of pre-processing techniques available, however there is not a clear consensus on an ideal pre-processing pipeline. Bad pre-processing can even cause predictive models to have a lower accuracy in comparison to no pre-processing^[80].

1.4.2.3 Comparison

Flow cytometry and Raman spectroscopy are two techniques that are able to give information about phenotypical traits of individuals in a bacterial population. In general, flow cytometry is a fast technique that gives information on cell-morphology and specific cell properties for which has been stained. Raman spectroscopy is a slower technique that gives a more holistic view on the molecular phenotypic traits. A short comparison of some of their main features is given in Table 1.2.

Table 1.2: Comparison of some of the main features of flow cytometry and Raman spectroscopy.

	Flow cytometry	Raman spectroscopy
Information	morphological features, information regarding nucleic acids content, cell viability, etc. (depending on the stain)	chemical composition of the cell, presence of biomolecules
Throughput	very high, up to 50,000 cells per second	much lower
Level of automation	use of autosamplers, possibility for online measurements	not very automated, requires more manual labor
Staining	mostly, not in case of autofluorescent cells	unnecessary

1.5 In silico communities

Following up the community composition over time is often a point of interest in synthetic ecology studies. But as mentioned in Section 1.2.2, this is not easily done using the classical sequencing or molecular techniques. Previous research revealed the potential of using machine learning techniques to distinguish between different bacterial^[81;82] or phytoplanktonic^[83] species based on their flow cytometric fingerprints, in both the artificial^[81] and the natural^[84] context.

Recently, a new method has been developed to assess community composition based on flow cytometric fingerprinting^[58]. In contrast to some of the previous studies, this method does not require any adaptations to the flow cytometer. The pipeline is illustrated in Figure 1.10. First, a fingerprint

of the axenic cultures that make up the synthetic community is made. In the next step the data of the axenic cultures is aggregated to a so-called ‘in silico community’. This in silico community consists of labelled data, which allows the use of supervised machine learning techniques. A classifier is trained to learn the difference between the fingerprints of the community-members. The label to be predicted is the species, while the features or predictors are the scatter and fluorescence parameters. Once this classifier has been trained on the dataset, it can be used to predict the relative abundances of the species in a mixture^[58]. As was the case in previous studies, some species are more difficult to be distinguished than others because of overlapping phenotypical traits.

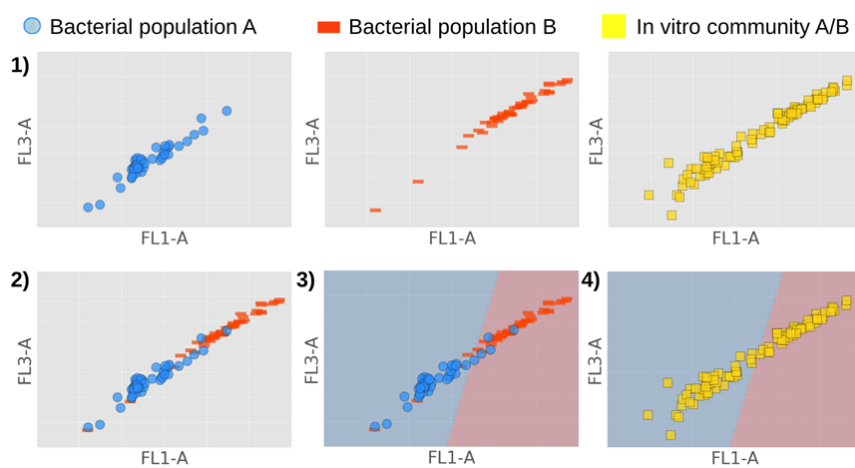


Figure 1.10: Conceptual illustration of a machine learning approach to infer the community composition of synthetic microbial communities. The first step is to acquire cytometric fingerprints of the axenic cultures that make up the synthetic community, as well as a cytometric fingerprint of the synthetic community under study. In the second step the community is created ‘in silico’ by aggregating the data of the axenic communities. In the next step a classifier is trained to learn the difference between the two species (in this example LDA). Finally, this classifier can be used to predict the abundances of the two species based on the cytometric fingerprint of the mixed community^[58].

1.6 Objectives

Recently, there has been a growing interest in understanding the ecological importance of phenotypes and phenotypic diversity in microbial communities. To a large extent it is still unknown which factors are influencing it and what the potential importance of this fine scale diversity might be in both natural and engineered ecosystems. In literature, most studies are either focusing on single species experimental set-ups or on entire communities. Typically these studies evaluate metabolic potential or respiration in function of imposed stressors.

Sympatric bacterial populations are populations that are present in the same environment and therefore encounter each other on a regular basis, which implies they have the potential to interact. The main objective of this study is to assess the influence of binary interactions between sympatric bacterial populations in function of the phenotypic diversity of the interacting organisms. Ecological theories will be tested using synthetic ecosystems, in which drinking water isolates will be used as model organisms. During this study the phenotypic diversity will be assessed using two single-cell techniques: flow cytometry and Raman spectroscopy. Both techniques measure widely different optical properties of individual bacterial cells. Complementarity or similarities between the two techniques will be evaluated. Currently there exist no validated computational approaches for characterizing phenotypic diversity by means of Raman spectroscopy. We will evaluate the potential to create such pipelines.

Quantification of organisms in synthetic ecology experimental set-ups is currently a difficulty. The available tools require either large sample volumes (sequencing) or the development of specific primers (qPCR). Recently, a new method has been developed to assess community composition based on flow cytometric fingerprinting. The last goal of this study is to apply and evaluate the performance of this methodology.

Materials and methods

2.1 Bacterial strains

Drinking water isolates were provided by Pidpa (Provinciale en Intercommunale Drinkwatermaatschappij der Provincie Antwerpen, Belgium). The isolates were grown in liquid minimal medium, M9 with 200 mg/L glucose as a carbon source.

2.2 Molecular analysis

To identify the drinking water isolates, DNA extraction was performed using the FastPrep DNA-extraction protocol^[85]. The DNA was amplified by polymerase chain reaction (PCR) on the 16S rRNA gene, with 63F and 1378R primers and the following thermal cycle: 5 min at 94°C, 30x (1 min at 95°C, 1 min at 53°C, 2 min at 72°C) and 10 min at 72°C. The resulting PCR products were purified and sent out for Sanger sequencing (LGC Genomics GmbH, Germany).

2.3 Microbial analysis

2.3.1 Bacterial growth

The growth curve for a bacterium in a certain set of conditions (T, growing medium, etc.) can be determined by inoculating a fixed number of cells in fresh medium and measuring the change in turbidity, and hence cell density, over regular time-intervals. This can easily be done using a spectrophotometer^[86]. The curve that results from this type of experiment has a typical sigmoidal shape, to which several models can be fitted. The logistic growth model is widely used and given

by Equation 2.1, where A is the carrying capacity, μ is the maximum specific growth rate and λ is the lag phase^[87].

$$y(t) = \frac{A}{1 + \exp\left(\frac{4\mu}{A}(\lambda-t)+2\right)}. \quad (2.1)$$

Growth curves for some of the drinking water isolates were determined. First, the bacteria were plated on nutrient agar (Oxoid, UK) plates. The next day, a single colony was picked and transferred to liquid minimal medium (M9 with 200 mg/L glucose as a carbon source). Two days later, cell densities in the liquid cultures were determined by flow cytometry and the cultures were diluted to a density of 2×10^6 cells/mL in sterile, 0.22 μm -filtered PBS. Double concentrated minimal medium was prepared. A 96-well plate was filled by adding 100 μL of double-concentrated medium and 100 μL of the diluted culture to each well. Each strain was prepared in triplicate. Blanks were prepared by adding 100 μL of PBS to the double-concentrated medium. Optical density (OD) at 600 nm was measured using a Tecan Infinite[®] M200 PRO multiwell plate reader (Tecan Trading AG, Switzerland) at time intervals of 15 minutes over a total period of 120 hours. The temperature was set at 28°C.

2.3.2 Flow cytometry

For flow cytometric analysis, the samples were diluted and stained with nucleic acid stains. The stains that were used are SYBR[®] Green I (SG, 100x concentrate in 0.22 μm -filtered dimethyl sulfoxide) for total cell analysis and SYBR[®] Green I combined with propidium iodide (SGPI, 100x concentrate in 0.22 μm -filtered dimethyl sulfoxide) for live-dead analysis. Staining was performed as described previously^[63], with incubation for 20 min at 37°C in the dark. Samples were analysed immediately after incubation. Two flow cytometers were used in this study.

A C6 Accuri[™] flow cytometer (BD Biosciences, Belgium), which was equipped with four fluorescence detectors (530/30 nm, 585/40 nm, >670 nm and 675/25 nm), two scatter detectors and a 20-mW 488-nm laser. The flow cytometer was operated with Milli-Q (MerckMillipore, Belgium) as sheath fluid.

A FACSVerse[™] flow cytometer (BD Biosciences, Belgium) with nine fluorescence detectors (527/32 nm, 783/56 nm, 488/15 nm, 586/42 nm, 700/54 nm, 660/10 nm, 783/56 nm, 528/45 nm and 488/45 nm), a scatter detector and a blue 20-mW 488-nm laser, a red 40-mW 640-nm laser and a violet 40-mW 405-nm laser. The flow cytometer was operated with FACSFlow[™] solution (BD Biosciences, Belgium) as sheath fluid.

2.3.3 Raman spectroscopy

The fixation protocol for Raman spectroscopy was adapted from a previously described protocol^[88]. 1mL of cell-suspension was centrifuged for 5 minutes at room temperature and 5000 g. The supernatant was discarded and the cell pellet was resuspended in cold, 0.22 µm-filtered PBS (4°C). The cell-suspension was again centrifuged for 5 minutes at room temperature and 5000 g. The supernatant was discarded and the cell pellet was resuspended in the fixative, 0.22 µm-filtered 4% (v/v) paraformaldehyde in PBS (pH 7.2). The sample was allowed to fix for 1h at room temperature, in the dark. The fixative was removed by centrifuging for 5 minutes at room temperature and 5000 g and resuspending the pellet in cold, 0.22 µm-filtered PBS (4°C), twice. The fixed sample was stored at 4°C. Prior to analysis, the fixed sample was centrifuged for 5 minutes at room temperature and 5000 g and the pellet was resuspended in 0.22 µm-filtered, milli-Q (4°C). The cell suspension was spotted onto a calciumfluoride slide and allowed to dry.

The dried sample was analysed immediately using a WITec Alpha300 R+ confocal Raman microscope with a 100x/0.9NA objective (Nikon, Japan), a 785 nm excitation diode laser (Toptica, Germany) and a UHTS 300 spectrometer with a -60°C cooled iDus 401 BR-DD CCD camera (Andor Technology Ltd, UK). Laser power before the objective was measured daily and was about 150 mW. Spectra were acquired in the range of 110-3375 cm⁻¹ with 300 grooves/mm diffraction grating. For each single cell spectrum, the Raman signal was acquired over 40 s.

2.4 Data analysis

2.4.1 Identification isolates

To identify the drinking water isolates, the sequences obtained via Sanger sequencing were blasted against two databases, NCBI blast and the Ribosomal database project. Both databases resulted in the same identities.

2.4.2 Bacterial growth

Prior to fitting the growth model, the optical density of the blanks was subtracted from each time series to correct for absorption caused by the medium. For each replicate the time series was selected from start to the point where the stationary phase was reached. Next, the data was imported

in R (v3.3.1)^[89] and analysed using the grofit package (v1.1.1-1)^[90]. The logistic growth model was fit to each replicate and R^2 for each fit was calculated to evaluate the goodness of fit.

2.4.3 Flow cytometry

Phenotypic diversity

The data was imported in R (v3.3.1)^[89] using the flowCore package (v1.40.3)^[91]. A quality control of the datasets was performed using the flowClean package (v.1.12.0)^[92]. An air bubble or large particle can potentially disturb the fluidics system for a very short or longer period of time and disturb the process of hydrodynamic focussing. The flowClean package contains an algorithm that tracks changes in fluorescence frequency within a sample. Events which are potentially anomalies are flagged and excluded for further analysis. After quality control, the background of the fingerprints was removed by manual gating on the primary fluorescent channels. A suitable gate was created for each experiment separately. Data was further analysed using the PhenoFlow package (v1.1)^[63]. This package translates the single-cell cytometric data into a phenotypic fingerprint and subsequently calculates established diversity metrics, including alpha- and beta-diversity as explained in Section 1.4.2.1. The Bray-Curtis dissimilarities were visualised by principal coordinate analysis (PCoA, for more information see Appendix 5.1.8). Statistical significance of the differences between the samples was assessed using permanova (for more information see Appendix 5.1.2).

In silico communities

After gating, the data was exported from R under Flow Cytometric Standard (FCS) format. The files were converted to comma-separated values files (CSV) to be further analysed in Python, using Scikit-learn^[93]. The classifier that was used for the in silico communities is a random forest (for more information see Appendix 5.1.4). To build the model, datasets were split into a balanced trainings- and testset (70%/30%). During building of the model, 200 trees were grown as this had previously been found to be sufficient^[58].

2.4.4 Raman spectroscopy

The data was analysed in R (v3.3.1)^[89]. Spectral preprocessing was performed using the package MALDIquant (v1.16)^[94].

To preprocess the Raman spectra, the biologically relevant part of the spectrum ($600\text{-}1800\text{ cm}^{-1}$) was selected^[76]. The spectra hold 333 data points over the selected range. Baseline correction was performed using the statistics-sensitive non-linear iterative peak-clipping (SNIP) algorithm^[95]. A high number of iterations was selected in order to make the result of the baseline correction less sensitive to small differences in the spectra. The spectra were normalised by means of surface normalisation, i.e. setting the surface under the spectrum equal to 1. Finally, the necessity for peak alignment was evaluated.

2.5 Experimental setups

2.5.1 Experiment 1

The aim of this experiment was to evaluate whether microbial interactions between sympatric bacterial populations can lead to changes in phenotype and phenotypic diversity of the interacting organisms. Drinking water isolates were used as model organisms for the sympatric bacterial populations. Two drinking water isolates were selected based on two criteria: the combination of these bacteria had good performance for the supervised in silico community methodology described in Section 1.5 and both bacteria were relatively fast growing (e.g. stationary phase was reached within 24 hours, starting from a cell density of 10^6 cells/mL). The selected bacterial species were *Enterobacter* sp. and *Pseudomonas* sp.. Since the organisms were merely used as model organisms they will further be denoted as taxon A and B.

During the experiment, bacteria were cultured in transwell plates (Corning[®] Costar[®] 6-well cell culture plates, Corning Incorporated) where apical and basal compartments were created using cell culture inserts (ThinCert[™] Cell Culture Inserts with pore diameter $0.4\text{ }\mu\text{m}$, Greiner Bio-One). The membrane of the culture inserts was replaced by membranes with smaller pore sizes to avoid migration of bacteria between the two phases (Cyclopore[®] polycarbonate and polyester membranes with $0.2\text{ }\mu\text{m}$ pore size, Whatman). Different combinations of apical and basal phase and corresponding starting cell densities are presented in Table 2.1. The starting cell densities were set to have an initial cell density of 10^6 cells/mL in each synthetic community.

Before the start of the experiment, both bacteria were plated on nutrient agar (Oxoid, UK) plates. From each plate a single colony was picked and transferred to liquid minimal medium (M9 with 200 mg/L glucose as carbon source). After two days of incubation at $28\text{ }^\circ\text{C}$, cell densities in the liquid cultures were determined by flow cytometry and diluted to the desired starting cell densities

in fresh medium. The required dilution was high enough to neglect differences in volume of fresh medium, and thus resources for growth, that were needed to prepare the cultures. Each combination of apical and basal phase was prepared in triplicate and randomised over the plates to account for plate effects. A blank, containing fresh medium in both apical and basal phase, was present on each plate as a control for cross-contamination. The plates were incubated at 28 °C and gently shaken (25 rpm) to aid diffusion of the metabolites between the compartments.

Table 2.1: Initial conditions in the 6-well plates. Cell densities were set for each synthetic community to have an initial cell density of 10^6 cells/mL and with equal relative abundances for both community members in the coculture and mixed culture.

Basal phase (4 mL)	Apical phase (2 mL)	Cell density basal phase (cells/mL)	Cell density apical phase (cells/mL)
A	fresh medium	1.5×10^6	0
A	B	7.5×10^5	1.5×10^6
B	fresh medium	1.5×10^6	0
A+B	fresh medium	$7.5 \times 10^5 + 7.5 \times 10^5$	0

The communities were monitored over a period of 72 hours. Every 24 hours samples were taken from each compartment for flow cytometric analysis. Each sample was split in two and diluted in 0.22 μ m-filtered, sterile PBS. One part was analysed using SG staining, the other part was stained with SGPI. The samples were analysed on the FACSVerse™ flow cytometer. After 72 hours samples for Raman spectroscopy were taken and fixated. All Raman samples were analysed within 1 week, with minimal time between them to limit possible differences caused by differences in duration of the storage.

2.5.2 Experiment 2

The aim of this experiment was to evaluate whether the influence of microbial interactions on phenotypic diversity, that was found in the previous experiment, was reversible. The same bacteria, and a similar set-up as for the first experiment were used.

Before the start of the experiment, both bacteria were plated on nutrient agar (Oxoid, UK) plates. From each plate a single colony was picked and transferred to liquid minimal medium (M9 with 200 mg/L glucose as carbon source). After two days of incubation at 28 °C, cell densities in the

liquid cultures were determined by flow cytometry and diluted to the desired starting cell densities in fresh medium. At the start of the experiment four synthetic ecosystems were created: two axenic cultures and two cocultures (Table 2.2). The axenic cultures were created in triplicate, the cocultures in sextuplicate. The communities were randomised over the plates to account for plate effects. After three days of incubation, the first flow cytometric measurement took place and the apical phases of the cocultures were replaced by new apical phases containing either 0.22 μm -filtered milli-Q or fresh minimal medium (Table 2.2). The plates were incubated at 28 °C. From the third day on, the cultures were monitored for another three days. Every 24h samples were taken for flow cytometric analysis. Each sample was split in two and diluted in 0.22 μm -filtered, sterile PBS. One part was analysed using SG staining, the other part was stained with SGPI. The samples were analysed on the FACSVerse™ flow cytometer.

Table 2.2: Initial conditions in the 6-well plates during the first experiment. Cell densities were set for each synthetic community to have an initial cell density of 10^6 cells/mL. The cocultures were created for both A and B in the basale phase since this was more practical for the replacement of apical phases on the third day. Two axenic cultures were present as reference for non-interacting genotypes.

Basal phase	Apical phase	Cell density basal phase (cells/mL)	Cell density apical phase (cells/mL)	Replacement (t = 72h)
A	fresh medium	1.5×10^6	0	none
B	fresh medium	1.5×10^6	0	none
A	B	7.5×10^5	1.5×10^6	milli-Q
A	B	7.5×10^5	1.5×10^6	fresh medium
B	A	7.5×10^5	1.5×10^6	milli-Q
B	A	7.5×10^5	1.5×10^6	fresh medium

2.5.3 Experiment 3

A *gfp*-labelled strain of the *Enterobacter* sp. used in experiment 1 was available. In order to check the evolution of the relative abundances that were found in the first experiment, mixed cultures of *Pseudomonas* sp. and the autofluorescent *Enterobacter* sp. were created. Before the start of the experiment, the bacteria were plated on nutrient agar (Oxoid, UK) plates. From each plate a single colony was picked and transferred to minimal medium (M9 with 200 mg/L glucose as carbon source). After two days of incubation at 28 °C, cell densities in the liquid cultures were

determined by flow cytometry. Three mixed cultures were created in 10 mL tubes, with equal initial cell densities as the mixed culture in the first experiment. The communities were incubated at 28 °C.

The mixed cultures were followed up over a period of 72 hours. Every 24 hours samples were taken for flow cytometric analysis. The samples were split in two and diluted in 0.22 µm-filtered bottle water (Evian). Due to the autofluorescent properties of *Enterobacter* sp. these cells can be detected on the primary fluorescent channel of the flow cytometer without staining. The first sample remained unstained and thus gives information about the cell density of *Enterobacter* sp. only. The second sample was stained with SG to get information on the total cell density. The samples were analysed on the C6 Accuri™ flow cytometer in fixed volume mode (25 µL).

2.5.4 Experiment 4

The aim of this experiment was to evaluate whether presence of multiple carbon sources would lead to a higher phenotypic diversity. Before the start of the experiment, *Enterobacter* sp. and *Pseudomonas* sp., that were used in previous experiments, were plated on nutrient agar (Oxoid, UK) plates. From each plate a single colony was picked and transferred to minimal medium (M9 with 200 mg/L glucose as carbon source). After two days of incubation at 28 °C, cell densities in the liquid cultures were determined by flow cytometry and the cultures were diluted to 2×10^6 cells/mL in sterile M9 without carbon source. Three different carbon sources were used (Table 2.3). For each of the carbon sources M9 with double carbon concentrations was prepared.

Synthetic communities were created in sterile 12 well-plates (Corning® Costar® 12-well cell culture plates, Corning Incorporated) by adding 2 mL of the diluted cultures to 2 mL of the double concentrated medium. The communities that were present on each of the plates are presented in Table 2.4, and were randomized over the plate to account for plate effects. The 12 well plates were incubated at 28°C. Every 24h samples were taken for flow cytometric analysis. Each sample was diluted in 0.22 µm-filtered bottle water (Evian) and stained with SG. The samples were analysed on the C6 Accuri™ flow cytometer using fixed volume mode (25 µL).

Table 2.3: Different carbon sources that were added to M9. For glucose a concentration of 200 mg/L was selected, similar to previous experiments. For the mixture of glucose, acetate and pyruvate the concentrations were determined to have a total C content equal to the one of the glucose treatment, and with equal contributions of the 3 compounds to the total C content. For yeast extract a concentration of 200 mg/L was selected. Since yeast extract has varying compositions and thus also varying C contents, the total C content in this treatment will likely differ from those of the other two treatments.

Carbon source(s)	Compound stock solution	Final concentration in M9 (mg/L)
Glucose	Glucose (C ₆ H ₁₂ O ₆)	200
Glucose, acetate, pyruvate	Glucose (C ₆ H ₁₂ O ₆)	66.67
	Sodiumacetate (CH ₃ COONa)	91.06
	Sodiumpyruvate (C ₃ H ₃ NaO)	81.44
Yeast extract	Yeast extract (Oxoid, UK)	200

Table 2.4: Initial conditions for the synthetic communities present on each of the three 12-well plates. On each plate all possible combinations of bacteria and carbon sources are present once. Each plate contains three wells with fresh medium as a control for cross-contamination.

Cell density A (cells/mL)	Cell density B (cells/mL)	Carbon source(s)
1 × 10 ⁶	0	glucose
0	1 × 10 ⁶	glucose
5 × 10 ⁵	5 × 10 ⁵	glucose
0	0	glucose
1 × 10 ⁶	0	glucose, acetate and pyruvate
0	1 × 10 ⁶	glucose, acetate and pyruvate
5 × 10 ⁵	5 × 10 ⁵	glucose, acetate and pyruvate
0	0	glucose, acetate and pyruvate
1 × 10 ⁶	0	yeast extract
0	1 × 10 ⁶	yeast extract
5 × 10 ⁵	5 × 10 ⁵	yeast extract
0	0	yeast extract

3.1 Experiment 1: Interactions between bacteria can lead to adjustment of their individual phenotypic diversities.

The aim of this experiment was to evaluate whether microbial interactions can lead to changes in phenotypic diversity of the interacting organisms and if this change is dependent on the composition of the sympatric bacterial populations. Phenotypic diversity was assessed through flow cytometry and Raman spectroscopy. Two drinking water isolates *Enterobacter* sp. and *Pseudomonas* sp. were selected as model organisms. Since the organisms were merely used as model organisms they will further be denoted as taxon A (*Enterobacter* sp.) and B (*Pseudomonas* sp.), respectively. Both isolates were relatively fast growing. Based on growth curves that were determined prior the start of the experiment, species A reached the stationary phase at 11h and species B at 22h.

Four synthetic communities were created (Figure 3.1). Both isolates were grown in axenic cultures as a reference for the non-interacting genotypes. To be able to study the individual community members separately after they have been interacting via their joint medium, a coculture with physical separation by a membrane was created. Lastly, a mixed culture without physical separation, representing ‘full interaction’, was created. Each community was created in triplicate. These synthetic communities were monitored over a period of 72h. Every 24h samples were taken for flow cytometric analysis and at 72h samples were taken for Raman spectroscopy. Since the first measurement took place at 24h, both community members were expected to be in stationary phase during the whole experiment. For flow cytometric analysis both SG and SGPI staining was applied. From the SGPI stained samples it can be concluded that all cultures were viable throughout the entire experiment. Following results are based on the SG stained samples (i.e. total cell analysis).

In the following results the samples are indicated with names in the form of ‘X treated with Y’, where X is the species in the sample (A, B or AB) and Y is what was present on the other side of the membrane (A, B or fresh medium). When referring to the ‘coculture’, the physically sepa-

rated mixed culture is intended, while ‘mixed culture’ indicates the mixed culture without physical separation.

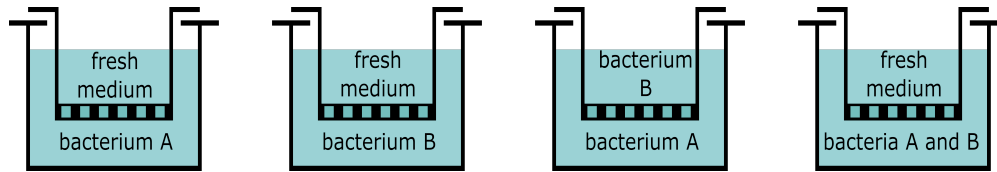


Figure 3.1: Illustration of the experimental set-up of experiment 1. Bacteria in apical and basal phase can interact via metabolites in their shared medium while they are physically separated by the membrane of the cell culture inserts. Four synthetic communities were created: two axenic cultures, a coculture and a mixed culture. For each synthetic community biological replicates ($n = 3$) were present.

3.1.1 Phenotypic diversity assessment through flow cytometry

The cell densities for each bacterial species as well as the mixed community were monitored at three time points. Cell densities of species A remain stable over time, while for species B and the mixed culture a limited increase over time is observed (Figure 3.2).

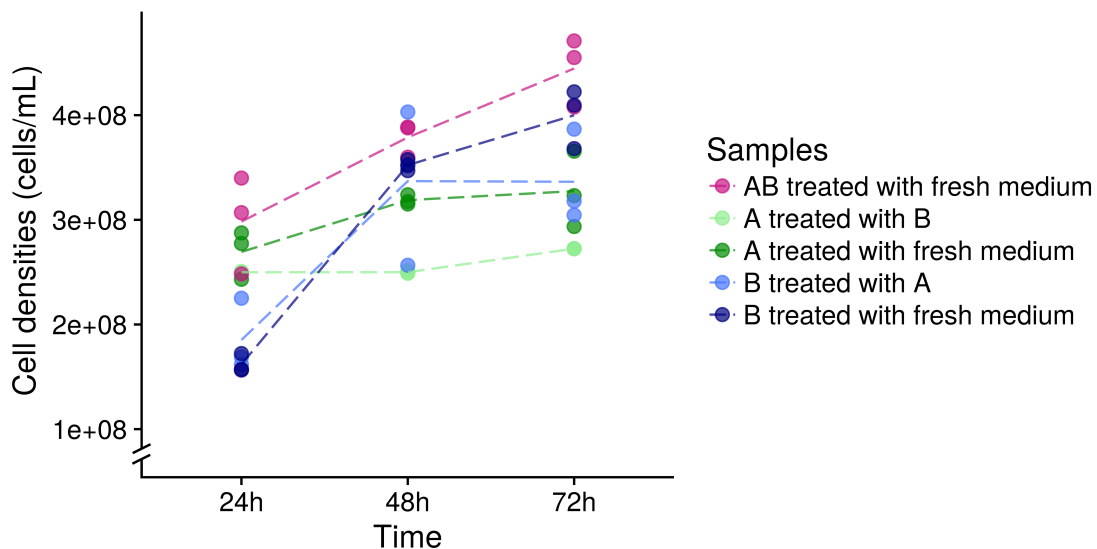


Figure 3.2: Evolution of cell densities for both individual bacterial species in communities of single species, cocultures and mixed cultures. There were biological replicates ($n = 3$) for each community. The dashed lines indicate the average trend of the replicates.

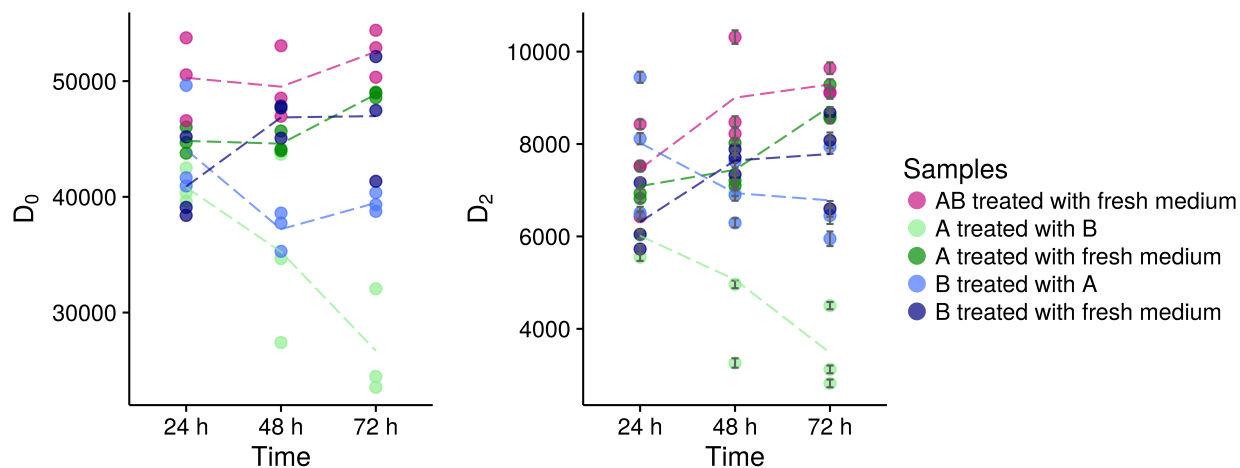


Figure 3.3: Hill diversity parameters D_0 (richness) and D_2 (richness and evenness) for both individual bacterial species in communities of single species, cocultures and mixed cultures. Error intervals on the D_2 are generated by bootstrapping (999 bootstraps). There were biological replicates ($n = 3$) for each community. The dashed lines indicate the average trend of the replicates.

Phenotypic alpha- (within sample) and beta- (between sample) diversity was calculated based on two scatter and five fluorescence detectors of the FACSVerse flow cytometer, as these were found to be the most informative^[96]. Results for D_1 are not shown since D_1 is highly correlated with D_2 ($r_p = 0.97$). When looking at phenotypic richness (D_0), the diversity of the species in axenic cultures is larger compared to the diversity of the same species when it was present in one of the compartments of the coculture (Figure 3.3). This difference is more pronounced for species A. The same observation is true for D_2 , with the exception of one replicate of B treated with A and one replicate of B treated with fresh medium. The difference between the axenic cultures and the coculture becomes larger over time. Diversity in the mixed community is generally largest, however, there are some exceptions and the difference is only small. The phenotypic richness of the mixed community remains relatively constant over time. It should be noted there is some variability in the diversity dynamics of the biological replicates.

To further compare the phenotypic fingerprints, a PCoA ordination was generated based on the Bray-Curtis dissimilarity between the fingerprints (Figure 3.4). Overall the populations are significantly differing ($p = 0.001$, $r^2 = 0.859$). Homogeneity of variance in groups of replicates was assessed before performing permanova. In this ordination, the fingerprints of both species are separated, with the mixed culture in between. The populations showed a significant shift in their phenotypic structure through time ($p = 0.001$, $r^2 = 0.158$). All populations evolve in the same direction over time (left to right in Figure 3.4), except for species A that was grown in axenic cul-

ture. In addition, there is a significant difference in the fingerprints of species A when present in an axenic culture compared to when present in the coculture ($p = 0.001$, $r^2 = 0.412$). For species B the differences in the fingerprints when present in an axenic culture compared to when present in the coculture were not significant ($p = 0.089$, $r^2 = 0.168$). The mixed culture shifted from a community that is more resembling A at the first measurement, towards a community that is more similar to species B at the second and third measurement. In general, the biological replicates are ordinated together, even when there was some variability in their alpha diversity dynamics.

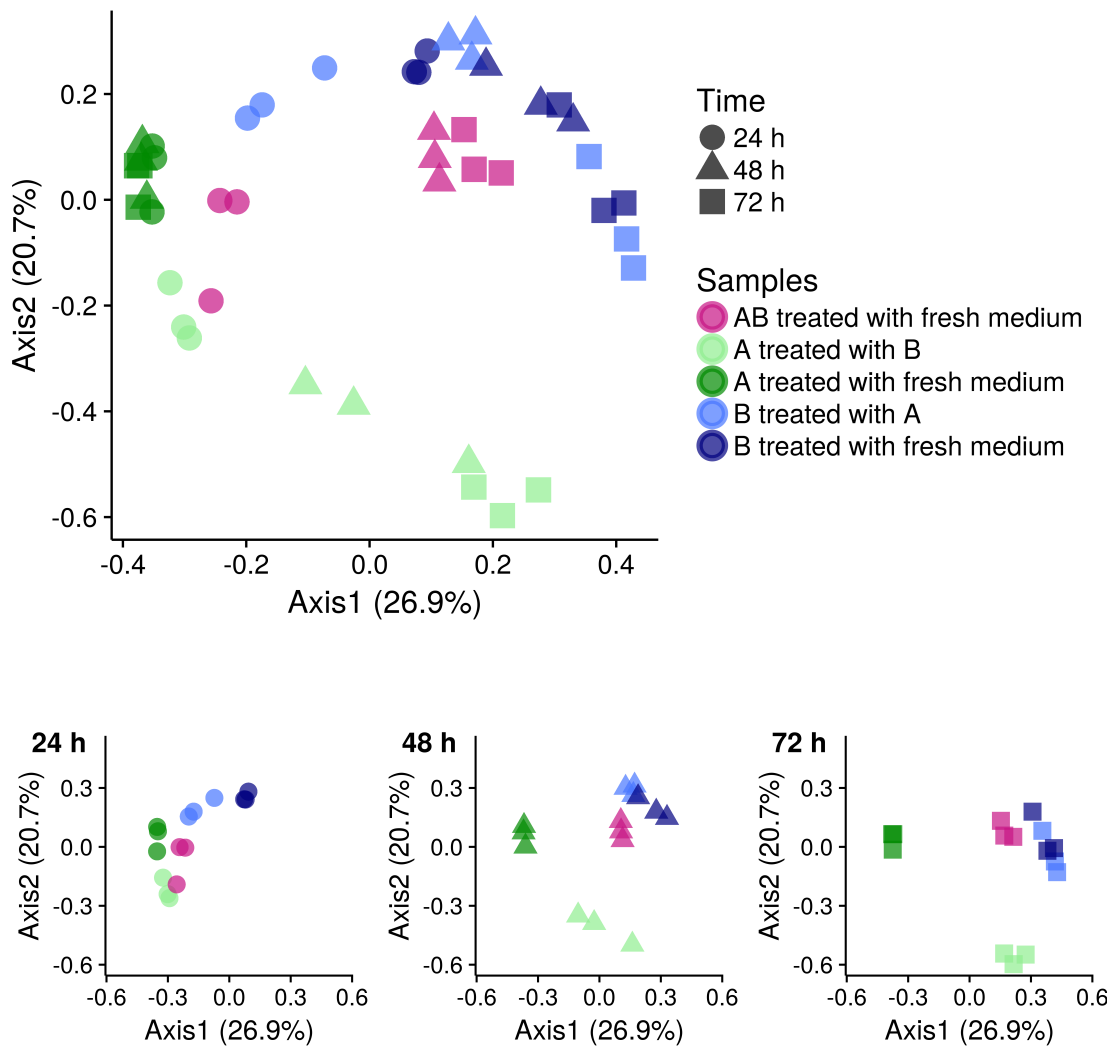


Figure 3.4: PCoA ordination of phenotypic fingerprints for all individual bacterial species in communities of single species, cocultures and mixed cultures. In the three lower graphs the result of the above ordination was split according to the different time points, since this allows for easier interpretation of how the different communities are relating to each other at each time point. There were biological replicates ($n = 3$) for each community.

To assess the shift in phenotypic community structure that is occurring due to the interaction, the differences in scattering pattern (FSC and SSC) and fluorescence (FITC) of the axenic cultures and the cocultures was evaluated using a contrast analysis. Differences in scatter patterns were limited for both species. In contrast, a clear difference in nucleic acid content was observed (Figure 3.5). For species A there is a shift towards high nucleic acid individuals in the coculture as compared to the axenic culture. This difference becomes larger over time. For species B there is a more limited difference, with a small enrichment of lower nucleic acid individuals. Thus, the shift in phenotypic community structure that is observed is depending on the taxon.

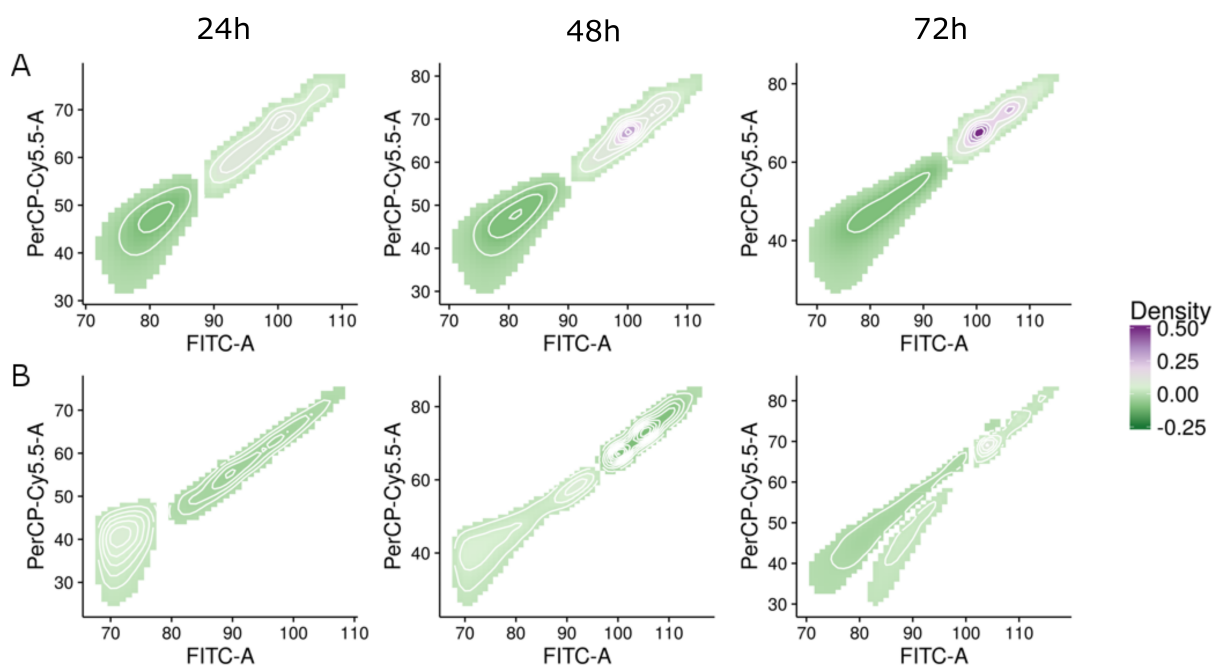


Figure 3.5: Contrast analysis of the phenotypic fingerprints to compare the difference in phenotypic community structure of axenic cultures and coculture members with respect to nucleic acid content. Each plot is a comparison between the axenic culture and coculture of the same species at the same time point, averaged over the biological triplicates. The color gradient indicates whether populations in the coculture increased (purple) or decreased relative (dark green) to their respective axenic growth at the specified timepoint. Pale green indicates no or very limited changes. The upper row (**A**) gives contrast results for species A, the lower row (**B**) gives contrast results for species B. If the difference between the two communities is lower than 0.01, no contrast value is shown on the graphs.

3.1.2 Prediction of relative abundances in the mixed culture

In the previous section the phenotypic structure of each taxon was analyzed when grown in axenic, coculture, and mixed culture conditions. In order to infer the community composition of the mixed cultures the supervised machine learning approach described in Section 1.5 was used. For training of the random forest, the biological replicates were pooled together and 10,000 cells of both A and B were randomly sampled. The data was partitioned into a balanced (i.e. the cell numbers for species A and species B are equal in these datasets) training and test set of 70% and 30% respectively. Relative abundances in the mixed culture were predicted for each of the biological replicates separately.

The analysis conducted above confirmed our initial hypothesis that the phenotypic diversity of a taxon can be conditional on the presence of other taxa. Therefore, in order to take into account this conditional aspect, we compared the abundance predictions of models that were constructed using different fingerprints as input data for model training (Figure 3.6, Appendix 5.3.1 for exact values). First, the random forest was trained on the fingerprints of the axenic cultures at the first measurement. Second, the random forest was trained, for each time point separately, on the fingerprints of the axenic cultures of the corresponding time point. Third, the random forest was trained, for each time point separately, on the fingerprints of the coculture members at the corresponding time point. The differences in predicted relative abundances between training the random forest on the data of the axenic cultures from the start or on the data of the axenic cultures from the corresponding time point, are small ($<1\%$). Both predict a higher abundance of A at 24h, and a community that hardly contains any member of A ($<1\%$) from 48h on. In contrast, the predictions based on the coculture members indicate a more gradual shift in community composition. As was the case for the predictions based on the axenic cultures, the predictions indicate a higher abundance of A in the community at 24h, but now a gradual enrichment of species B at the second and third time point is predicted. This indicates that the choice of input data leads to different predicted relative abundances.

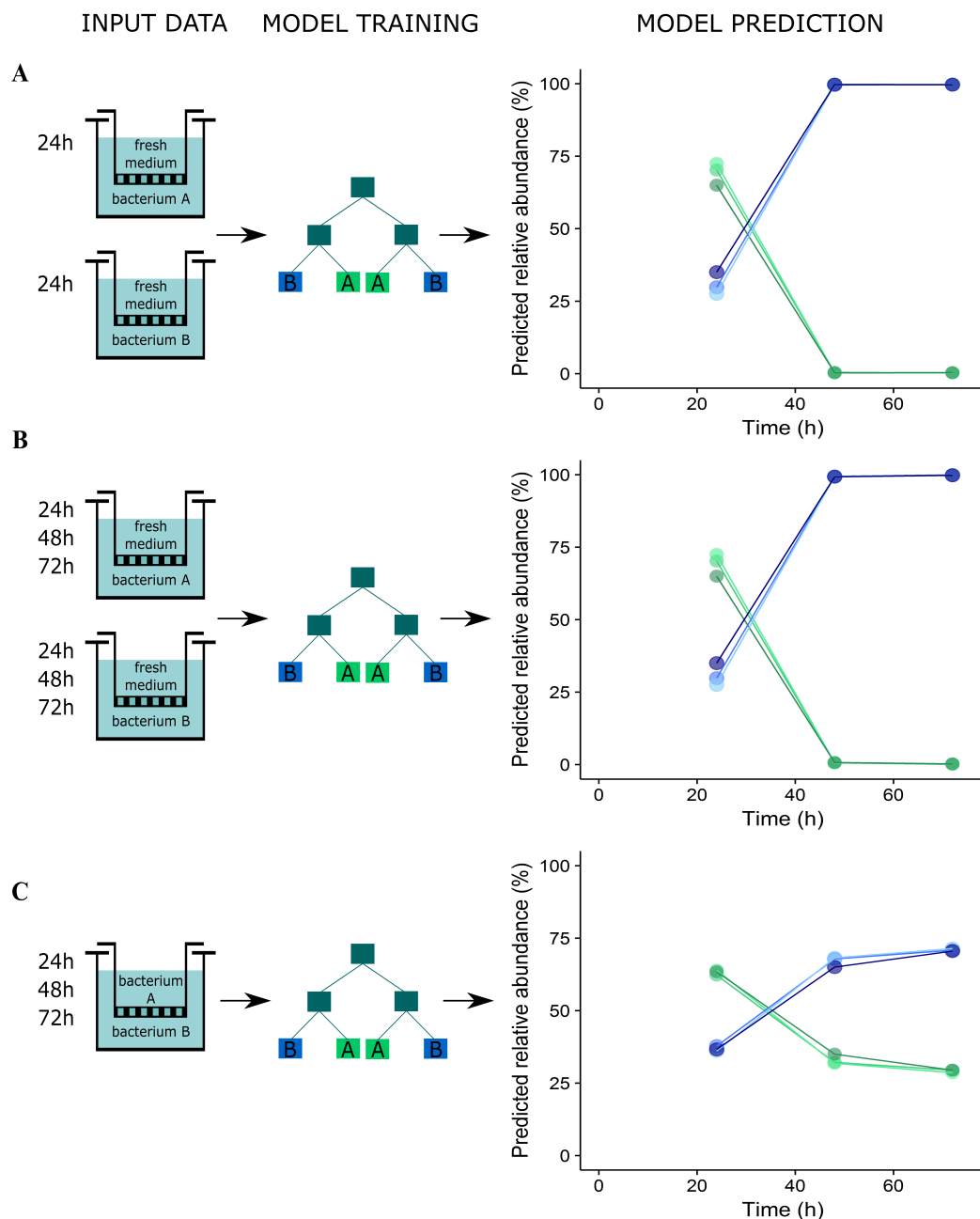


Figure 3.6: Relative abundances in the mixed cultures were predicted using the supervised in silico community methodology. The random forest classifier that was used to infer community composition of the mixed cultures ('AB treated with fresh') was trained using different input data, as represented on the left side. **A:** The random forest was trained using the fingerprints of the axenic cultures from the first measurement. **B:** The random forest was trained on the fingerprints of the axenic cultures of the corresponding time point. For example, the relative abundances in a sample at $t = 48h$ was predicted by a random forest that was trained using the data of the axenic cultures of A and B at $t = 48h$. **C:** The random forest was trained on the fingerprints of the coculture members at the corresponding time point.

3.1.3 In silico approach

In the previous section it was found that using data from different communities as input data for the model led to different predicted relative abundances. In order to evaluate which of the above predictions are most correct and whether the interaction in the coculture is different from the one in the mixed culture, the random forest was used to create communities on the computer (in silico) by merging or splitting data of the measured samples (Figure 3.7). Communities were created in two ways: data of the mixed cultures were split into ‘in silico coculture samples’ and data of either axenic cultures or cocultures were merged together to ‘in silico mixed culture samples’.

To create the mixed cultures in silico, data of the coculture members were mixed according to the relative abundances that were predicted by the random forest that was trained on members of the coculture. For example, to create an in silico community of the mixed culture at 48h, data from a measurement of species A in the coculture at 48h were concatenated to data of species B at 48h in the coculture in a ratio of 33/67, as these were the relative abundances that were predicted by the random forest that was trained on the cocultures at 48h (Figure 3.6, C). The same approach was taken to mix the data of the axenic cultures according to the predicted abundances of the random forest that was trained on the fingerprints of the axenic cultures (Figure 3.6, B).

To create the cocultures in silico, mixed communities were split in two. To decide whether a cell in data of the mixed culture belonged to species A or to species B a decision threshold was set for both species separately. These thresholds were determined based on the ROC curve of the random forest. The ROC curve (receiver operating characteristic) indicates the true positives and false positives of a binary classifier in function of the threshold that is used to discriminate between the two classes (for more information see Appendix 5.1.6). In case of the random forest, the standard approach is to take the majority vote among the decision trees, thus using a threshold of 0.5. Rather than selecting the threshold for which no wrong predictions were made (i.e. FP = 0 on the ROC curve), the threshold was chosen so that the corresponding point on the ROC curve was as close as possible to the point (0,1). This way the random forest was allowed to make some mistakes in classifying the cells. The reason for this choice is that if the fingerprints of both species have overlapping phenotypes, and the random forest is not allowed to make any mistakes, these overlapping phenotypes can not be classified as A, neither as B. This way a ‘hole’ would be created in the fingerprints of the in silico communities where the overlapping phenotypes occur. The false positive rate, i.e. wrong classifications, at the selected thresholds was <0.5% (Appendix 5.3.2). The random forest classifier was trained on cells of the pooled replicates, but

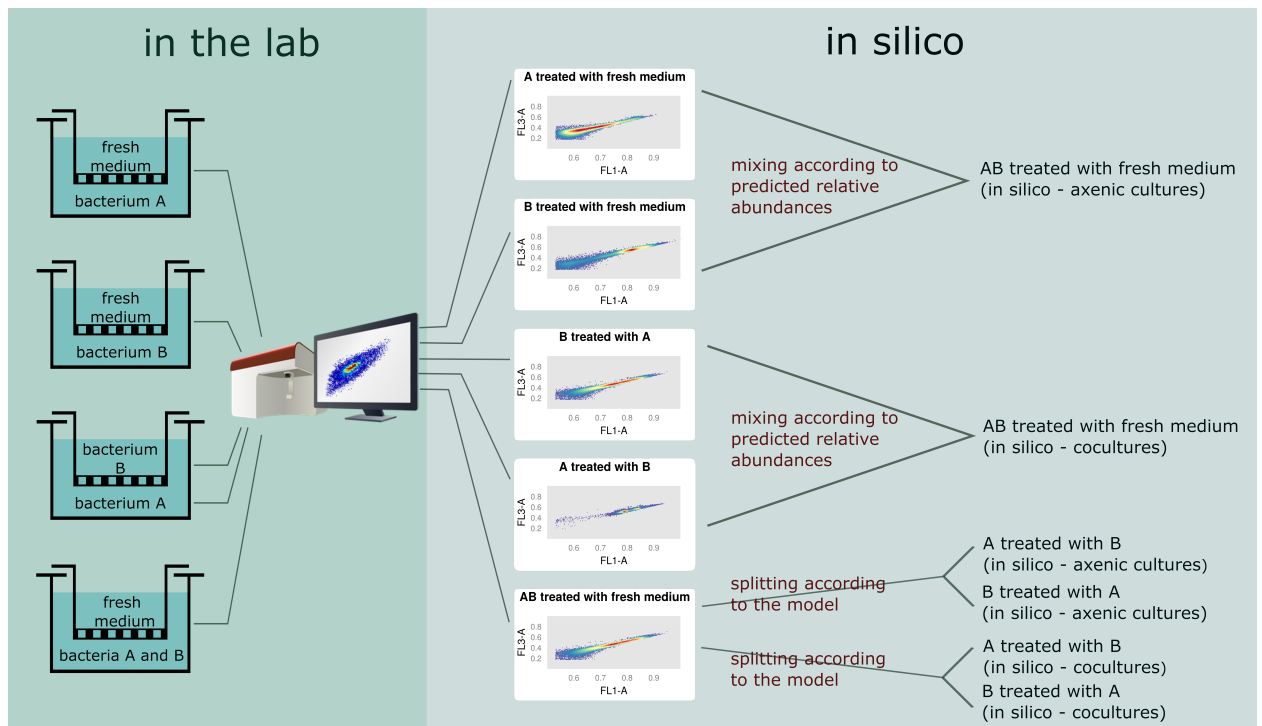


Figure 3.7: Illustration of the approach to create communities in silico. Data from FCM measurements is split or merged to create communities in silico. These communities were created in two ways: data of the mixed cultures was split into in silico coculture populations and data of either axenic cultures or cocultures was merged together to in silico mixed culture populations. The mixed cultures were created based on the relative abundances as predicted by the random forest. The cocultures were created by splitting the mixed samples according to optimal thresholds that were determined based on the ROC curve. Between brackets is indicated whether the community was created based on the model or predictions of the axenic cultures or the cocultures.

when creating the communities the biological replicates were not pooled. Also during merging of the data of the coculture members, the data of cultures that were present in the same well was merged, thus there was no combination of data from the apical phase of one well to data of the basal phase from another well. This way the in silico created communities can still be thought of as three biological replicates. The phenotypic fingerprint of these in silico communities could then be compared to the ones in the measured samples in order to evaluate whether the fingerprint of the species in the coculture is the same as the fingerprints of both members of the mixed cultures, i.e. whether interaction in the coculture and interaction in the mixed culture lead to the same changes in phenotypic fingerprint. This allowed to validate our experimental set-up.

Fingerprints of the in silico created communities were compared to those of the measured samples by procrustes analysis on the PCoA ordinations. In short, this technique compares two matrices

by translating, reflecting, rotating and dilating one of the matrices in order to minimize the sum of residuals between the two matrices. A randomized test is used to determine the significance of the obtained sum of residuals (see Appendix 5.3 for more details). This technique allows to quantify whether the ordination based on the relation between the *in silico* created communities is similar to the one of the measured samples (Figure 3.8). Since no counterparts for the axenic cultures can be created *in silico*, the original measurements of the axenic cultures was appended to the datasets of *in silico* cultures. This way also the relationship of the *in silico* created communities towards the axenic cultures was taken into account. Due to the low predicted relative abundances of species A when using the random forest that was trained on data of the axenic cultures, the ‘*in silico* - axenic cultures’ for species A have low cell numbers. During diversity estimation the samples of each population are randomly resampled to the lowest sample size, which led to communities of only 99 cells. Therefore the results of the procrustes analysis for the ‘*in silico* - axenic cultures’ and ‘*in silico* - cocultures’ cannot be compared directly. Both ordinations are showing a trend in time, similar to the one of the measured samples. The ordinations are respecting the relationships towards the axenic cultures of species A. The relationship towards the axenic cultures of B seems to be preserved better for the ‘*in silico* - cocultures’ compared to the ‘*in silico* - axenic cultures’. For the ‘*in silico* - cocultures’, the samples of the mixed cultures are no longer ordinated between the ones of the cocultures, they are ordinated in between the samples of species B. It should be noted that the procrustes analysis evaluates replicates of the same community as completely different instances. During the minimisation of the sum of residuals a match will be made between replicate 1 of the *in silico* samples for a certain community to replicate 1 of the measured samples. However, replicate 1 of the *in silico* samples is not more related to replicate 1 of the measured samples than would be replicate 2 or 3 of the *in silico* samples. Thus there is some artificial matching of samples. However, since replicates are ordinated close to each other and this issue is occurring in both the comparison of measured samples to ‘*in silico* - axenic culture’ samples and the comparison of measured samples to ‘*in silico* - coculture’ samples, this was not expected to be a problem. Both ordinations were found to be significantly similar to the one of the original measurements (Table 3.1).

Table 3.1: Procrustes analysis (999 permutations) was used to compare the ordination of the in silico created communities to the one of the measured samples. The m^2 statistic indicates the final concordance of both ordinations. The smaller m^2 , the better the fit. (* = ordinations based on only 99 cells, other ordinations are based on 20,000 cells)

Ordination 1	Ordination 2	m^2	p-value
Measured*	In silico - axenic cultures*	0.834	0.001
Measured	In silico - cocultures	0.885	0.009

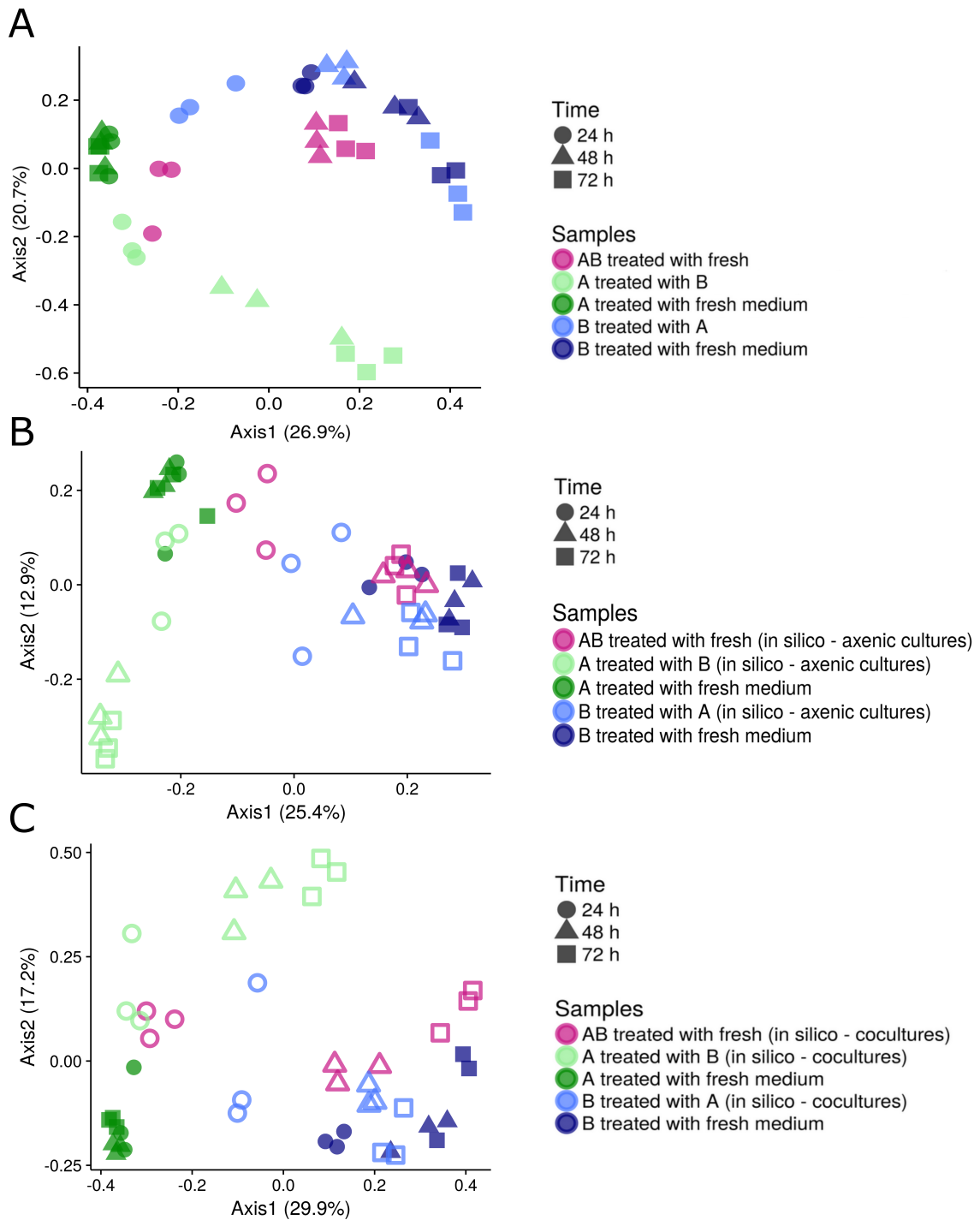


Figure 3.8: **A:** Ordination of the measured samples. **B:** Ordination of the measured axenic cultures and the in silico created samples that are based on the random forest that was trained on data of the axenic cultures. **C:** Ordination of the measured axenic cultures and the in silico created samples that are based on the random forest that was trained on data of the cocultures. Note: the in silico created samples that are based on the random forest that was trained on data of the axenic cultures (**C**) is based on only 99 cells, other ordinations are based on 20,000 cells. (● = measured, ○ = created in silico).

3.1.4 Phenotypic characterization through Raman spectroscopy

Through FCM analysis it was found that the phenotypic community structure of a taxon is depending on the community composition. Raman spectroscopy was used to measure single cell spectra for phenotypic characterization. Raman spectra of single cells for the axenic cultures and the two coculture compartments were acquired from one of the replicates of each community at 72h in the experiment. For each sample between 51 and 55 single cell spectra were measured. To have a similar level of uncertainty, 51 spectra of each sample were selected for further analysis. The spectra with the lowest intensity were assumed to be of lesser quality, therefore the spectra with the lowest average intensity were discarded. After preprocessing, the average spectrum and corresponding standard deviations for each sample were calculated (Figure 3.9). A large peak in the range of $810 - 1010 \text{ cm}^{-1}$ was present in the average spectrum of species A in the axenic culture, while this peak was not observed in any of the other samples. The cells of this sample showed large differences in intensity of this peak, as can be seen by the large standard deviations. We believe this peak might be the result of technical issues during fixation or storage of the sample, and is not directly related to presence of biomolecules in the sample.

To gain insight in the differences between the phenotypic properties of cells from each microcosm, spectra were centered and scaled to perform PCA (Figure 3.10 A). Because of the presence of the large peak in the sample of species A in the axenic culture, the difference in these spectra compared to the other spectra seems to be large overall. PCA did not result in a clear separation of the samples. There is a large overlap between cells from species B that were grown in axenic culture and cells from species B that were grown in coculture, and a more limited overlap between the cells of species A that were grown in coculture and the cells of species B that was grown in axenic culture. Therefore another visualisation technique, t-SNE, was used to gain more insight in the underlying structure of the data (Figure 3.10 B). t-SNE is a highly performant visualisation technique that can map the underlying structure of high dimensional data in a low dimensional (2- or 3D) space. It specifically aims to preserve the 'local' structure of the data. By seeding the algorithm with the result of a PCA the 'global' structure can be respected as well (for more information see Appendix 5.1.9). Using t-SNE cells from each microcosm can be separated better. The spectra of both species are separated well and the spectra of the same species with a different treatment are separated relatively well too, with some exceptions for species B.

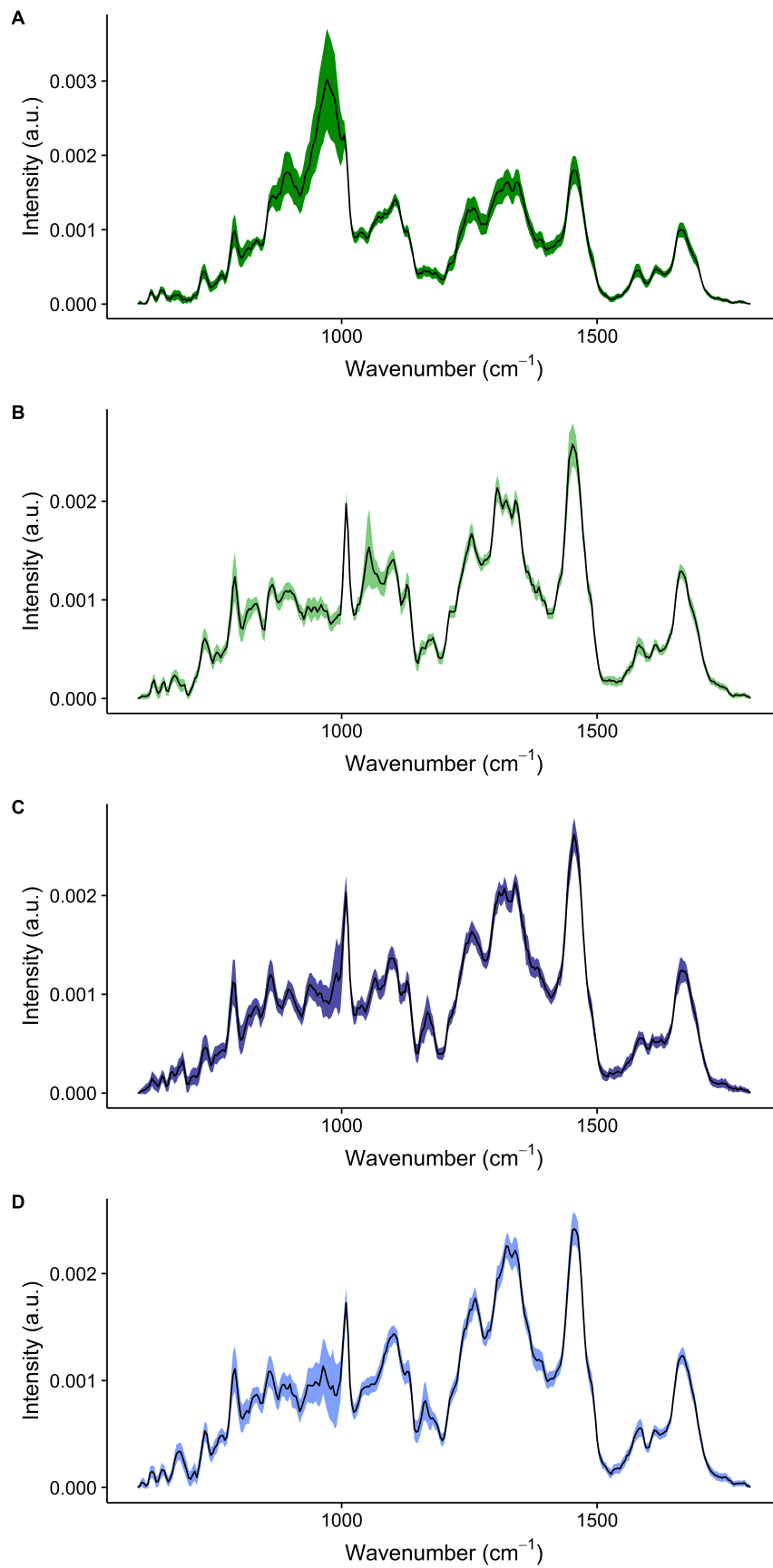


Figure 3.9: Average Raman spectra of the single-cell measurements. **A:** Species A in axenic culture. **B:** Species A in the coculture. **C:** Species B in axenic culture. **D:** Species B in the coculture. Colored bands indicate the standard deviations. All average spectra are based on 51 single cell measurements.

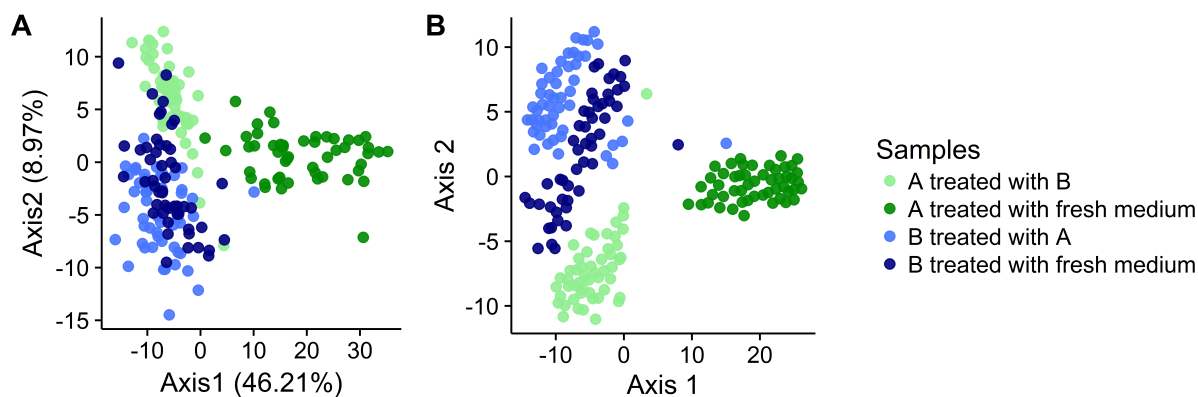


Figure 3.10: Visualisation of the underlying structure in the dataset of single cell Raman spectra for species A in axenic culture, species A in the coculture, species B in axenic culture and species B in the coculture. 51 single cell measurements are present for each of the samples. Visualisation was carried out using two techniques, PCA (**A**) and t-SNE seeded with the result of PCA (**B**). Note: Distances between clusters on a t-SNE plot cannot be interpreted as a measure of dissimilarity between the clusters.

3.1.5 Difference in cellular composition

In the previous sections, it was found that interaction between bacteria was leading to a shift in nucleic acid content and that the single cell spectra of the different populations can be separated. However, based on FCM results alone it is impossible to infer whether this shift in nucleic acids is related to a change in DNA or RNA content of the cells. To know which biomolecules were differing in spectra of the different samples, the location of peaks in the spectra can be compared to known databases of biomolecules, provided that these were measured using a laser with the same wavelength as the laser that was used to measure the single cell spectra. A database containing 60 preprocessed spectra of biological molecules measured with a laser with equal wavelength as the one used in this study is available for public from a study of De Gelder *et al.*, 2007^[77]. This database includes the spectra of all DNA and RNA bases (Appendix 5.2.1). Since the standard deviations around the average spectra indicated some variability between cells of the same sample, a robust method to find the wavenumbers to discriminate between the spectra of different treatments is needed. Finding these most important features to discriminate between classes is called ‘feature selection’. Several techniques to do this exist. Since we believe the peak in the range of 810 - 1010 cm^{-1} that was observed in the spectra of species A in the axenic culture is the result of technical issues, we excluded this region during analysis of the spectra of species A.

First, a randomized logistic regression model (RLR) was built for each species separately. The features are the intensity values for each of the wavenumbers and the label to be predicted is the treatment of the bacterium (axenic culture or coculture). A randomized model generates a ranking of features according to their importance for predicting the label (for more information on this procedure see Appendix 5.1.5). Based on this ranking several models were built. First, a model containing only the most important predictor was built. Next, a model containing the two most important predictors was built. This procedure was repeated until all predictors are used in the model. For each of the models a 10-fold cross-validation was carried out. The model that holds only the best predictor, will likely have too less information to make a reliable distinction between the classes and thus will not have the best accuracy. When adding other important predictors, more information is being incorporated in the model and thus the accuracy will go up. The point where the highest attainable accuracy is reached can be determined based on the cross-validation errors. This point corresponds to the optimal minimum number of required predictors. To evaluate whether these results could be generalized, a leave-one-out cross-validation (LOOCV) was applied. In LOOCV one sample of the dataset is set aside before building the model. When the RLR model was built and cross-validated, the optimal model was used to predict the label of the hold-out sample. Since this sample was not used during feature selection and validation, the accuracy on the hold-out sample gives information on whether the feature selection that has been found to be optimal was generalizable to unseen data. The entire procedure for feature selection and cross-validation was repeated for each of spectra as the hold-out sample. Since we believed the range of 810 - 1010 cm^{-1} in the spectra of species A in the axenic culture was due to technical issues and not biological differences, this region was removed in all spectra of species A. For both species only a small number of features was required to reliably predict their treatment (Table 3.2). Moreover, the selection of the low number of optimal features could be generalized to unseen data as can be seen by the high LOOCV accuracies.

Table 3.2: Results of the RLR models for both species. The cross-validated accuracies were determined using 10 folds.

Species	Range of optimal number of features	Cross-validated accuracy at optimal number of features (%)	LOOCV accuracy (%)
A	1 - 2	100	99.02
B	4 - 6	100	99.02

The RLR model showed that a clear distinction between treatments can be found for both species. However, since the spectra of biomolecules such as nucleic acids have several peaks over the biologically relevant region, the low amount of selected wavenumbers made it difficult to attribute the differences to certain compounds. Therefore another approach was taken, the χ^2 statistic was calculated for different treatments of the spectra for each species (Figure 3.11). The values for χ^2 were then visually compared to the spectra of the nucleic acid bases (Appendix 5.2.1). For most peaks in the spectra of the nucleic acid bases, both DNA and RNA, a peak in the values for χ^2 is found. Next to this, there were also peaks in the χ^2 -values that were not corresponding to peaks in of nucleic acid bases, which might indicate that also other compounds next to nucleic acid bases were differing between the treatments. However, one should be aware of the fact that each peak in a Raman spectrum is the result of multiple compounds, and thus from this analysis it is still somewhat ambiguous to conclude which biomolecules were differing between the treatments.

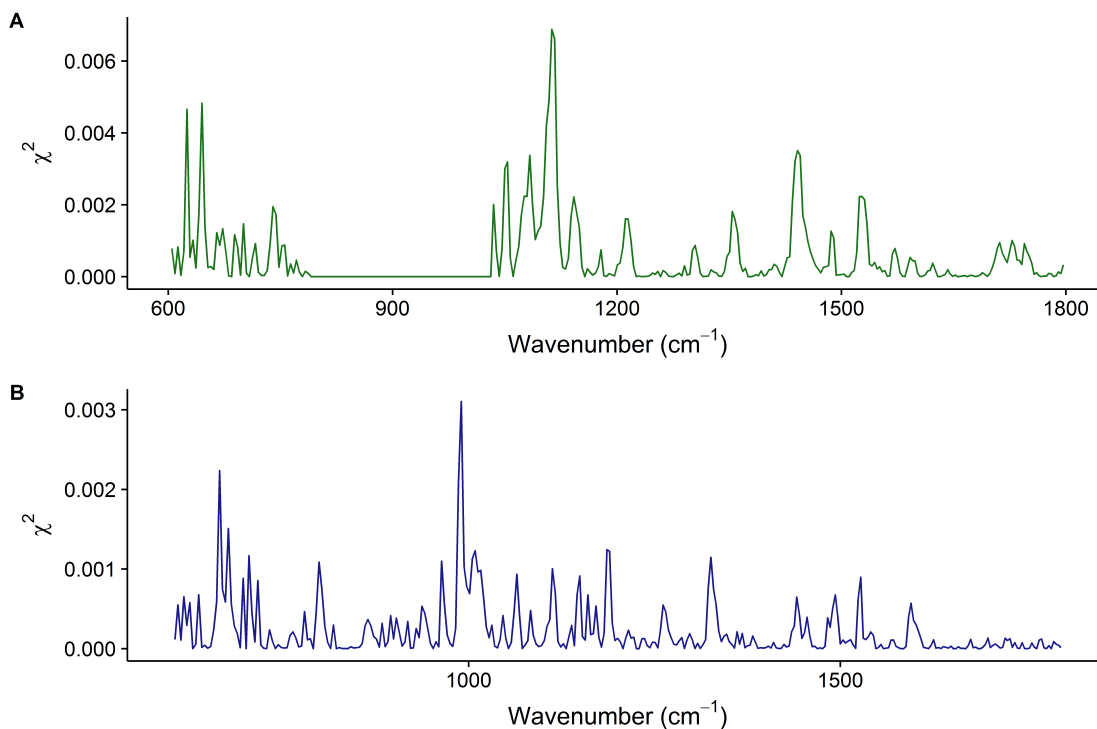


Figure 3.11: The χ^2 statistic between the spectra of the different treatments for each species separately. **A:** χ^2 statistic for species A. **B:** χ^2 statistic for species B.

3.2 Experiment 2: Reversibility of the effect of interactions on the individual phenotypic diversities of the bacteria.

The aim of this experiment was to evaluate whether the influence of microbial interactions on phenotypic diversity, that was found in the previous experiment, was reversible. The same bacteria, and a similar set-up as for the first experiment were used (Figure 3.12). Axenic cultures were created as a reference for the level and dynamics of diversity in non-interacting genotypes. Cocultures were created in the same way as during the first experiment. After three days of incubation, the apical phases of the cocultures were replaced with new apical phases, containing either fresh medium or milli-Q. Milli-Q was used since the level of dilution has been found to influence community dynamics, and thus only removing the apical phase and not replacing it by a new apical phase might lead to community dynamics that are related to the level of dilution or upconcentration and not to the presence or absence of an interacting partner^[97]. Fresh medium was used since the bacteria are already in stationary phase, and therefore might need some pulse to start recovering. After three days of incubation the first FCM measurement took place and the apical phases were replaced. The communities were then followed up over three days. Every 24h samples were taken for FCM and analysed using both SG and SGPI staining. From the SGPI stained samples it can be concluded that all cultures were viable throughout the entire experiment. Following results are based on the SG stained samples.

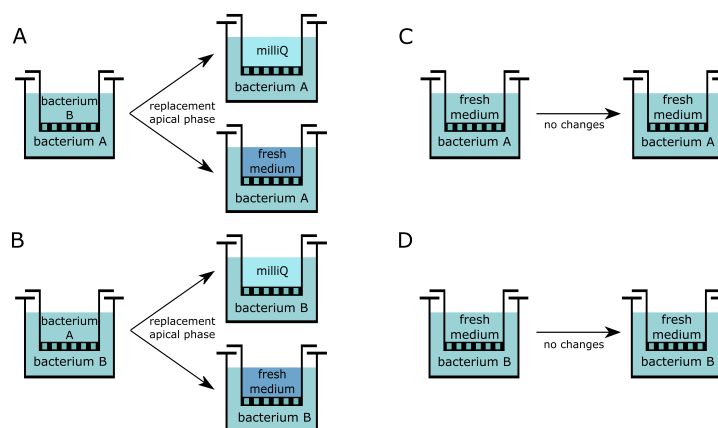


Figure 3.12: Illustration of the experimental set-up of the second experiment and the related hypothesis. Cocultures were created in the same way as for the first experiment. After three days of incubation, the apical phases of the cocultures were replaced with new apical phases, containing either fresh medium or milli-Q (A and B). Axenic cultures were created as a reference for non-interacting genotypes. For the axenic cultures there was no replacement of the apical phase (C and D).

Diversity was again calculated based on two scatter and five fluorescence detectors of the FACS-Verse flow cytometer^[96]. After three days, the expected lower level of phenotypic richness D_0 and phenotypic diversity D_2 was not observed for either of the species (Figure 3.13). Diversity of the coculture members was similar to the diversity of the axenic cultures. These diversity values are in the same range as the axenic cultures during the first experiment. Over time, a decrease of diversity in the axenic cultures ('control') is observed for both species. For species A, the replacement of the apical phase with either fresh medium or milli-Q led to an increase in phenotypic diversity over time. For species B, an increase in diversity was observed for the communities with fresh medium in apical phase, but not for the communities with milli-Q.

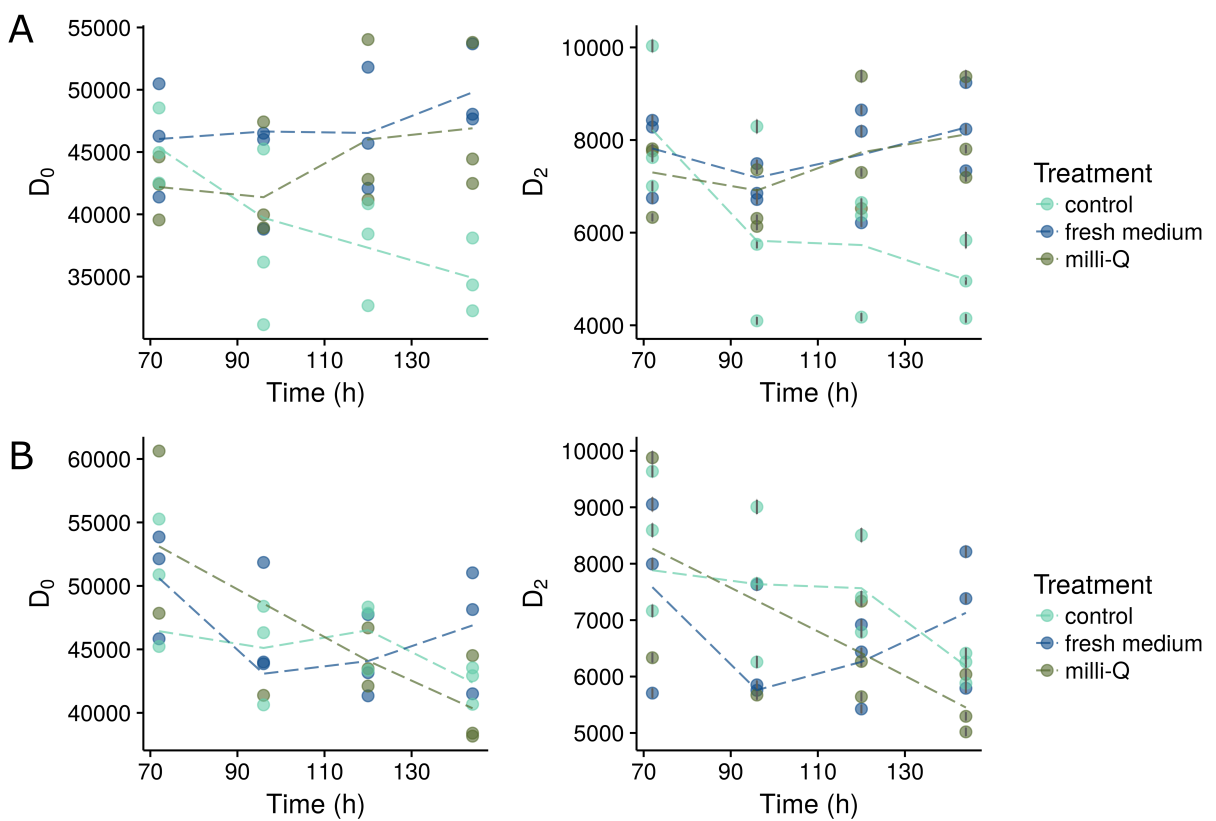


Figure 3.13: Hill diversity parameters D_0 (richness) and D_2 (richness and evenness) for both individual bacterial species in communities of single species, denoted as ‘control’, and cocultures, denoted as ‘fresh medium’ and ‘milli-Q’, from the third day on. The colors indicate the contents of the apical phase that was applied to the communities on the third day. The dashed lines indicate the average trend of the replicates. The upper row (A) gives results for species A, the lower row (B) gives results for species B. Note that the measurements started only on the third day

3.3 Experiment 3: Validation of the predicted relative abundances.

The aim of this experiment was to validate the trend in relative abundances that was found in the first experiment using the supervised in silico community methodology. A *gfp*-labeled strain of the *Enterobacter* sp. used in the first experiment was available. Mixed cultures of *Pseudomonas* sp., denoted as B, and the *gfp*-expressing *Enterobacter* sp., denoted as A, were created. These cultures were sampled every 24h over a total period of 72h, similar to the first experiment. Due to the autofluorescence properties of species A these cells can be detected on the primary fluorescent channel of the flow cytometer without staining. The samples were analyzed once without staining to get information about the cell density of species A, and once with SG staining to get information on the total cell density. This way the relative abundances of the two species in the mixed culture can be assessed. It should be noted that the *gfp*-labeled strain suffers from a metabolic burden due to the presence of his *gfp* and that not all bacteria with a *gfp*-label will be fluorescent, which might cause deviations as compared to what the relative abundances of the non fluorescent *Enterobacter* sp. were during the first experiment. Over time a gradual enrichment of species B was observed (Figure 3.14). This trend is most similar to what was found in the first experiment by training the random forest on the coculture members, however the absolute values of the relative abundances are differing. The predicted abundance of A was lower than what was found in this experiment.

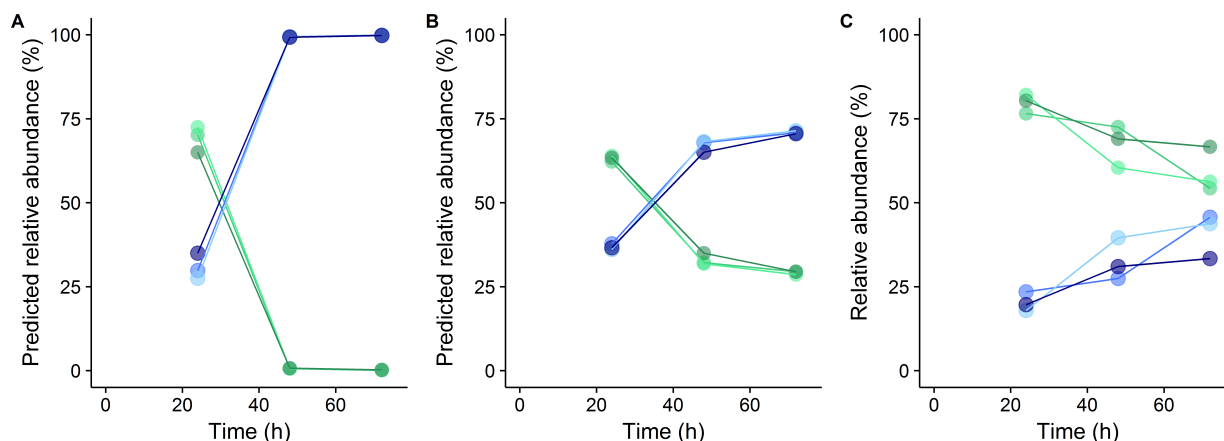


Figure 3.14: Evolution of the relative abundances as predicted by the random forest that was trained on fingerprints of the axenic cultures in experiment 1 (A), on fingerprints of the coculture members in experiment 1 (B) and as determined using a fluorescent A in experiment 3 (C). Green lines indicate the relative abundances of species A, blue lines indicate the relative abundances of species B. The different shades correspond to biological replicates ($n = 3$).

3.4 Experiment 4: Influence of carbon source diversity on phenotypic diversity.

The aim of this experiment was to evaluate whether presence of multiple carbon sources would steer the microbial populations to a higher phenotypic diversity. The hypothesis for this experiment was that the presence of multiple carbon sources leads more niche differentiation within the community, i.e. the breakdown and consumption of different carbon sources compared to only one, and therefore would lead to a higher phenotypic diversity. Three carbon sources were added to minimal medium: glucose, a mixture of glucose, acetate and pyruvate and yeast extract. For glucose and the mixture of glucose, acetate and pyruvate, the concentrations were set to have the same amount of carbon in each of the treatments. Since the composition of yeast extract differs among batches, the amount of C in the yeast extract treatment was not controlled. The same *Enterobacter* sp. and *Pseudomonas* sp., denoted as A and B, as for the first and second experiment were used. The two axenic cultures and a coculture were monitored over a total period of 72h. Approximately every 24h samples were taken for flow cytometric analysis.

Diversity was calculated based on two scatter and two fluorescence detectors of the Accuri C6 flow cytometer. Both D_0 and D_2 were evaluated, since it was unknown if presence of multiple carbon sources would lead to a higher level of phenotypic richness (D_0), a reorganization of the phenotypic community landscape (D_2), or both. Results for D_1 are not shown since D_1 is highly correlated with D_2 ($r_p = 0.99$).

When evaluating differences in diversity, no clear influence of the presence of multiple carbon sources was found (Figure 3.15). For both axenic cultures and the coculture, replicates of the different treatments have similar diversity values. When evaluating differences in diversity for each carbon source separately, similar diversity dynamics are observed for each species, irrespective of the carbon source (Figure 3.16). Except for yeast extract, where the diversity of species B at the first measurement is slightly higher compared to the diversity of species B in the other two carbon treatments. Thus no clear dependence of phenotypic diversity on the carbon source diversity was observed in this experiment.

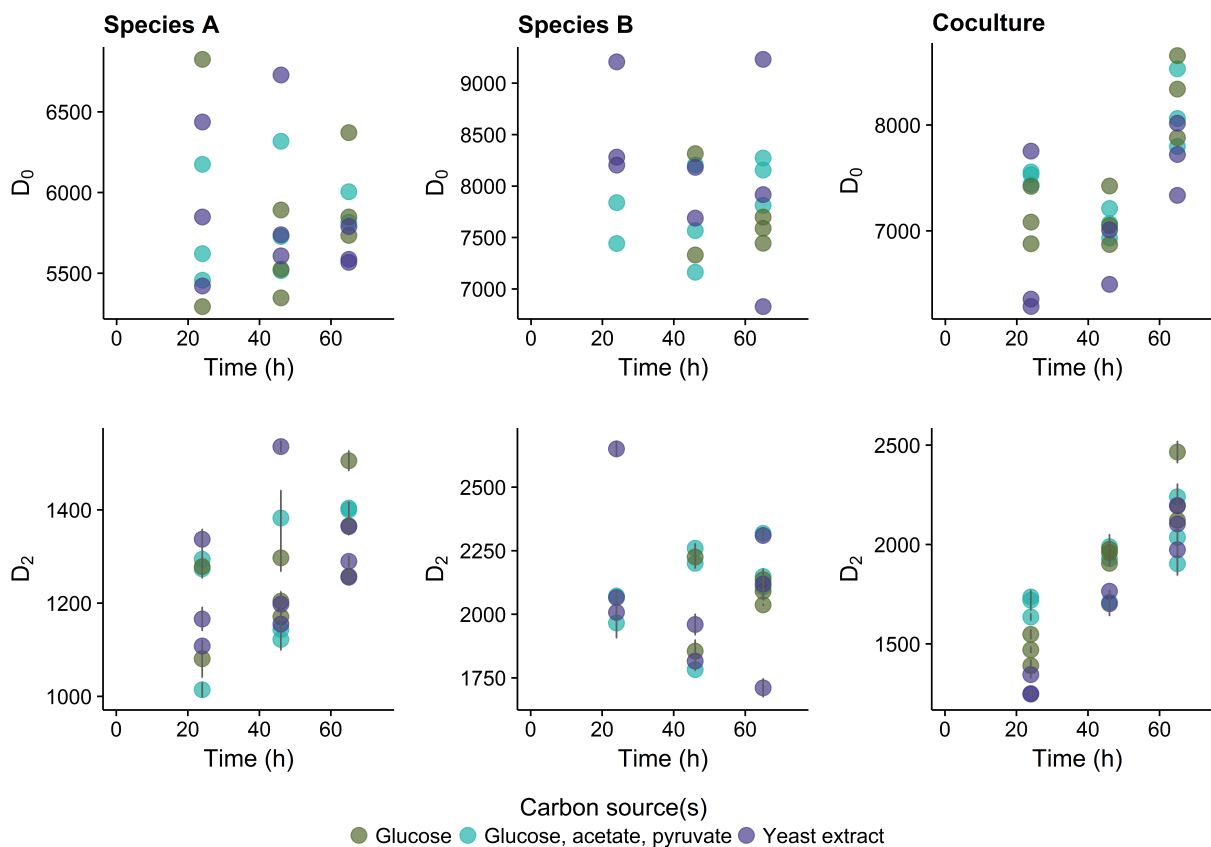


Figure 3.15: Hill diversity parameters D_0 and D_2 of both individual bacterial species and the mixed culture, for each carbon source separately. Error intervals on the D_2 are generated by bootstrapping (999 bootstraps). There were biological replicates ($n = 3$) for each sample, missing replicates are due to too low cell numbers for reliable diversity estimation. Note: y-axis are not set equal, since this made interpretation of the graphs more difficult.

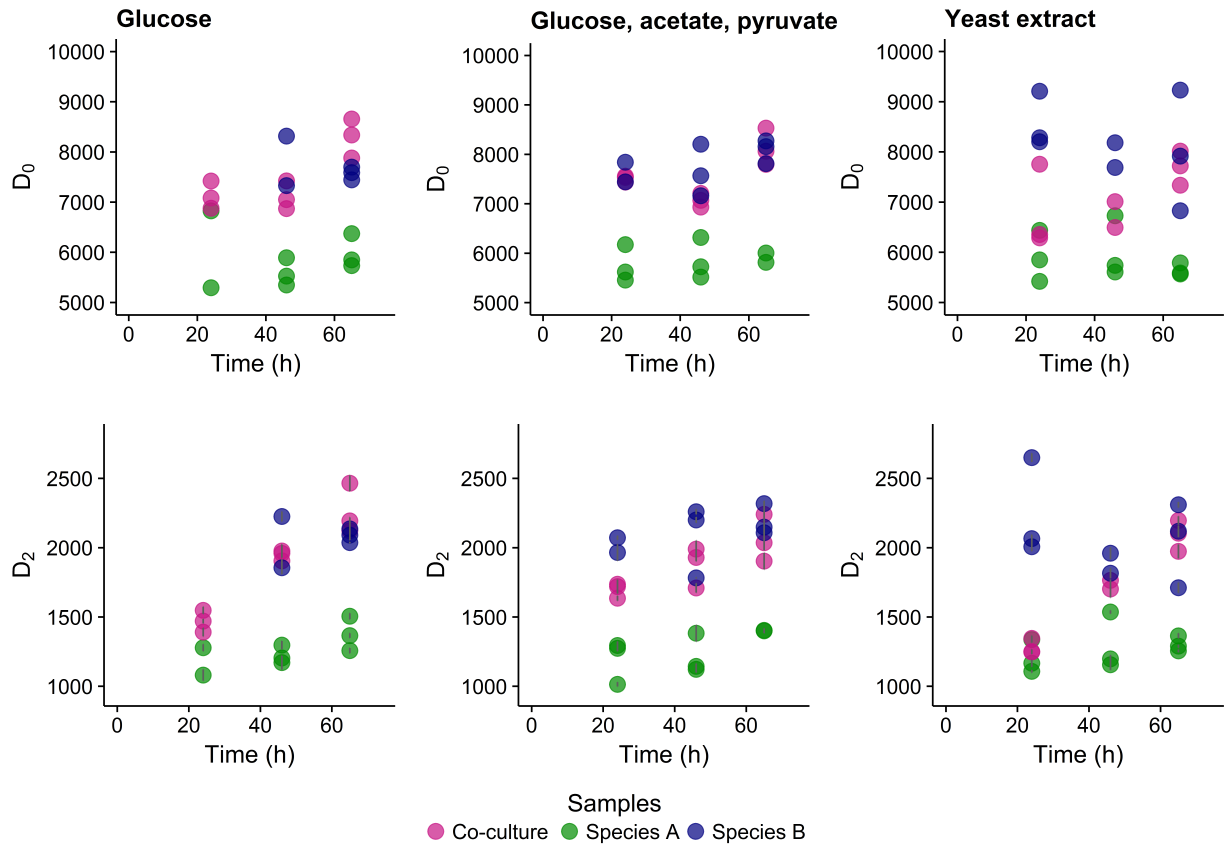


Figure 3.16: Hill diversity parameters D_0 and D_2 of both individual bacterial species and the mixed culture, for each carbon source separately. Error intervals on the D_2 are generated by bootstrapping (999 bootstraps). There were biological replicates ($n = 3$) for each sample, missing replicates are due to too low cellnumbers for reliable diversity estimation.

4.1 Hypothesis 1: Interactions between bacteria leads to an adjustment of their individual phenotypic diversities.

One of the aims in this study was to evaluate whether microbial interactions can lead to changes in the individual phenotypic diversity of the interacting organisms. To investigate this, an experiment was carried out where four synthetic communities were created, using two drinking water isolates as model organisms (Section 3.1). The communities that were created consist of two axenic cultures, a coculture with physical separation between the genotypes and a mixed culture without physical separation. In this experiment phenotypic diversity was evaluated through flow cytometry and Raman spectroscopy.

Based on the FCM results of this first experiment, the phenotypic diversity of the community members was different when they were grown in a coculture compared to when they were grown in axenic cultures. For both genotypes a lower phenotypic diversity was observed in the coculture compared to the same genotype in axenic culture (Figure 3.3). This effect of interaction on alpha diversity was more pronounced for species A than for species B, indicating that different genotypes have different phenotypic responses to the interaction. The phenotypic richness (D_0) decreased, which indicates that the interaction did not only lead to a reorganization of the phenotypic community structure (i.e. change in the relative abundances of the phenotypes), but that there are phenotypes which disappeared due to the interaction. Over time, the differences in alpha diversity between the axenic and cocultures was becoming more pronounced. This might indicate that the interaction which causes the adaptation of the individual phenotypic diversities is mainly present during stationary phase, and less pronounced or absent during the exponential growth of the bacteria.

When evaluating the beta-diversity between the populations, similar trends were found (Figure 3.4). The axenic cultures of species A were showing only moderate changes over time, while species A in the cocultures was more dynamic and became more distinct from the axenic cultures over time. For species B, the fingerprint of the coculture and the axenic culture followed a similar trend over time, with insignificant differences between the axenic culture and the coculture. This confirms our previous findings that different genotypes adopt different phenotypic responses to the interaction and that the interaction which causes the adaptation of the individual phenotypic diversities mainly takes place during stationary phase. However, for species B the difference between the two conditions was slightly more pronounced at 24h. This could be explained by the fact that species B was expected to reach stationary phase only just before the first measurement, and the moment at which species B reached stationary phase in the axenic culture and in the coculture might have been different. If the growth dynamics in the two microcosms were different, then the phenotypic fingerprints would also be different since different phenotypes are exhibited during the different stages of growth^[45]. The preliminary determined growth curves were based on OD measurements. Since OD measurements are not accounting for changes in cell properties, there is often a discrepancy between growth curve estimations based on OD measurements and the ones based on total viable counts^[98]. Therefore the exact point where the stationary phase was reached might be slightly different from our preliminary estimate and species B might not have been in stationary phase at 24h. For species B and the mixed culture of A and B a limited increase in cell density over time was observed, especially from the first measurement to the second (Figure 3.2). This might imply that species B was indeed not yet in stationary phase at 24h and thus differences in growth dynamics might be the explanation for the larger difference between the cocultures and the axenic cultures at 24h.

To further investigate the differences in phenotypic community structure based on the FCM results, scattering and fluorescence patterns of the populations were evaluated for each species separately (Figure 3.5). The differences in scattering patterns were limited for both species. Since scattering gives information on size, morphology and granularity^[52], our findings indicate there were no large changes in cell morphology due to the interaction for either of the species. The fluorescence signals are the result of the SG staining. SG is a nucleic acid stain that primarily stains double stranded DNA, but will also stain the RNA^[99]. Since SG staining is a stoichiometric staining^[50], a higher fluorescence signal is directly related to a higher nucleic acid content. In terms of nucleic acid content, large differences were observed for species A and smaller differences for species B. Species A is shifting to a community with a lower variance in nucleic acid content and with a higher

abundance of high nucleic acid individuals. This higher nucleic acid content might be the result of different shifts in physiology that could have been occurring. On the one hand, the cells could have had a higher DNA copy number, indicating an adaptation of their cell cycle. Bacteria are known to adapt their cell cycle behaviour and chromosome content under certain environmental conditions. For example, for *Pseudomonas putida* a constant growth rate with an accelerated DNA replication has been observed in relation to different types of stress^[100]. On the other hand, the bacteria might have had a higher RNA content, indicating a shift in their gene expression. The bacteria could have been more active while still expressing the same genes as they were in the axenic cultures, or they might have shifted towards expression of other genes compared to the axenic cultures. The spread in nucleic acids content in a clonal population can be attributed to several causes. There is natural stochasticity in gene expression, which is called ‘gene expression noise’^[101]. This noise can be attributed to noise in the expression itself and in other cellular components that are present in low concentrations. This biological noise has been found to be both controlled^[102] and structured^[103;104], i.e. the amount of biological noise in gene- and protein expression is assumed to be related to the function of the gene or protein. The noise levels are found to reflect the potential costs and benefits related to the noise, for example, stress-related proteins are more noisy compared to proteins for synthesis^[103]. Thus, if the bacteria were changing their gene expression due to the interaction, the noise levels of their gene expression might have changed as well. Based on flow cytometry results alone, it is impossible to draw conclusions on the reason for this difference in nucleic acid contents.

Through Raman spectroscopy we investigated whether specific biomolecules could be associated with the observed shifts in phenotypic diversity. The spectra of all DNA and RNA bases were available from literature and were measured using the same laser wavelength as in this study^[77]. We aimed to investigate which of the above mentioned scenario’s was most likely to be occurring. If only an intensity shift in the peak regions related to the DNA bases was observed, this would indicate there is only a change in the DNA content, implying an adaptation of the cell cycle. Uracil is replacing thymine in RNA molecules compared to DNA^[105]. If a shift in intensity related to the peaks in the spectrum of uracil was observed, a shift in presence of RNA and thus gene expression is indicated. In addition, this shift in gene expression would then also cause a difference in protein content, which would be detectable through Raman spectroscopy as well. It is not possible to draw a conclusion on whether this shift in gene expression is related to only a more intense activity or a different activity compared to the coculture. There is variation in the spectra of bacteria from the same population, indicating we can already observe a range of phenotypes

via the spectra. Therefore, we attempted to find the major differences in the spectra between different populations through feature selection (Section 3.1.5). Differences in intensities in the peakregions of both DNA and RNA bases were observed, as well as in regions where DNA and RNA bases have no peaks. These last might be related to the backbone of the DNA and RNA molecules or to other cell constituents such as proteins, etc. Since we observed changes in DNA bases as well as in RNA bases, we were still not able to conclude which of the above mentioned hypothesis was most plausible. Through more specific techniques such as transcriptomics, which would allow comparison of gene expression patterns, a better understanding of the interaction that was occurring might be obtained.

The phenotype of an organism is related to its functionality^[34]. Making the link between individuals and their activity is a difficult task and is one of the key objectives in microbial research nowadays^[106]. Based on the FCM fingerprints and Raman spectra we were not able to make this link either. However, if the assumption that a higher cytometric diversity (e.g. diversity in nucleic acid content and morphology) corresponds to the capacity to occupy a broader range of niches can be withhold, we can state an alternative hypothesis based on the results of this experiment: when a single genotype is growing in an axenic culture it is creating an entire ‘community’ of isogenic cells, where the constituent isogenic cells occupy multiple niches, leading to a high phenotypic diversity. In case multiple genotypes are present, two in this study, the available niches can be distributed among the genotypes according to their efficiency to occupy the niche. This would enable each genotype to occupy the niches at which it is most performant, thus creating a mixed community with a high functionality. This hypothesis of phenotypic niche differentiation due to genotypic richness is illustrated in Figure 4.1. During this study, communities containing only two genotypes were used. This means that no conclusions can be drawn regarding the validity of this hypothesis or the exact phenotypic diversity trajectory related to the amount of genotypes in the community. The observed change in diversity was different for both model organisms, indicating that diversity dynamics might be different for each genotype.

Besides the number of genotypes, there might be other factors that could potentially influence the diversity dynamics of a single genotype in the community. For example, the change in phenotypic diversity might be interaction-, partner- or environment-dependent. A community of species A and species B might lead to a different change in phenotypic diversity for species A, compared to when the community consists of species A and another species C. The environment might influence the amount of available niches, for example by its physical organisation and the existence of micro-environments^[107]. The amount of available niches for a certain genotype might also be a function

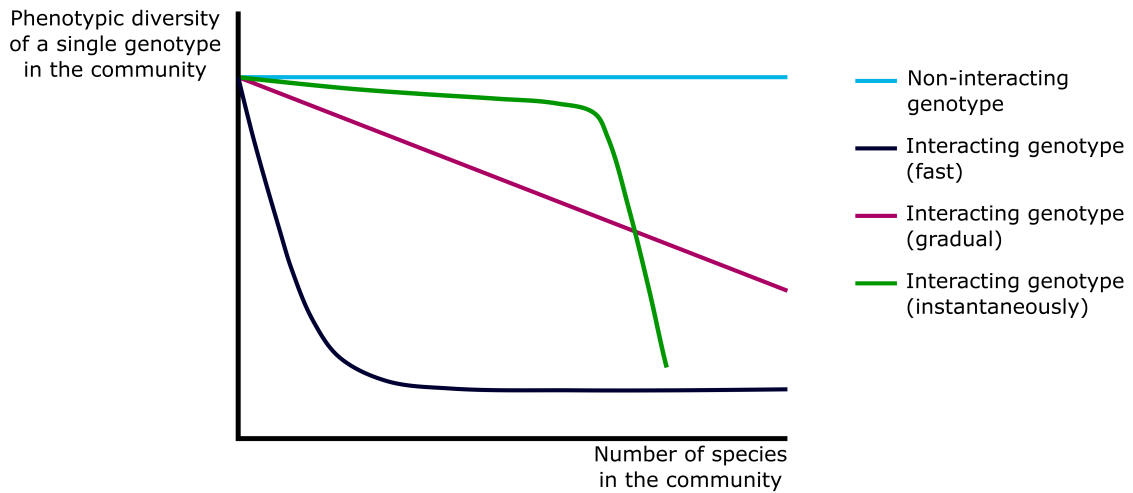


Figure 4.1: Potential trajectories for the hypothesis of phenotypic niche differentiation for a single genotype as a function of genotypic richness of the community. A non-interacting genotype might be unaffected by the presence of other species (light blue line). An interacting species might evolve to a lower phenotypic diversity where it is only occupying the niches where it is most efficient. This could happen at a fast rate (dark blue line), more gradually (pink line) or instantaneously when some other species that preferentially occupies similar niches enters the community (green line). Note that these are some hypothetical dynamics and that there might be other possible patterns as well.

of which other genotypes are present. For example, a species that is excreting a metabolite might make new niches available to species which can feed on this metabolite and in this way increase the phenotypic diversity of these species. Alternatively, a species which degrades antibiotics, can create antibiotics-free micro-environments and therefore create more available niches for non-antibiotic resistant organisms. Potentially, when a predator enters the community, the competition between various species may be controlled, possibly creating additional niches for some species.

Niche differentiation in sympatric populations is a phenomenon that has widely been observed for macro-ecological communities. It has been reported for two^[108] and more^[109] sympatric populations. Some patterns regarding diversity and functionality relationships in micro-ecology have been found to show similarities to those in macro-ecology^[110]. Niche differentiation might be a pattern that also occurs in microbial communities. This has previously been hypothesized in many studies, for example in a study regarding co-occurrence of sympatric yeasts related to their potential for resource partitioning^[111] and in a study where the relationship between gene expression similarity and the potential for interaction and co-existence of freshwater green algae was evaluated^[112]. Surprisingly, in this last study the researchers found that algae with more similar general gene expression patterns were more likely to co-exist while the initial hypothesis stated the op-

posite. Therefore they stated that their hypothesis might have been incomplete and hypothesized that the potential of species to differentiate their niches and to coexist is related to differences in expression of rare genes, rather than to the expression of the 'core genome', which is responsible for survival and reproduction. Since FCM and Raman spectroscopy both result in general measurements of nucleic acids, we can only hypothesize about potential differences and shifts in gene expression that were occurring here. As stated previously, more sensitive techniques such as transcriptomics might provide a better understanding of the community dynamics that were observed during this study.

Reversibility of the effect of interactions on phenotypic diversity

Related to this first experiment, two more experiments were carried out (Section 3.2 and 3.4). In the second experiment we attempted to assess the potential of the species to partially return to their axenic phenotypic diversity states after they were disconnected from their partner genotype. This hypothesis is illustrated in Figure 4.2. When the species are isolated from the coculture the microbial interaction ceases and they might return to the situation where they were creating an entire 'community' of isogenic cells, displaying a higher level of phenotypic diversity. Since the expected lower diversity for the species when grown in coculture was not observed after three days, we were not able to investigate this (Figure 3.13). The reason why we did not observe this lower level of diversity might be that there was a small difference in the experimental set-up of the first experiment compared to the second experiment. In the first experiment the cultures were gently shaken to aid diffusion of metabolites between the two compartments. In the second experiment there was no shaking, which might have led to a reduction of metabolite exchange, and thus the two species might have been unaware of each others presence. After three days, the apical phases were replaced with either milli-Q or fresh medium. This replacement will have caused some liquid exchange between the compartments which explains why some trends were observed after this replacement, even though the cultures were still not shaken.

For species A, higher levels of phenotypic diversity were observed in both treatments compared to the control cultures where the apical phase had not been replaced (Figure 3.13). This might be explained by the fact that the cultures were in stationary phase at high cell densities (3×10^8 cells/mL) and therefore the community might have been shifting towards lesser active states, exhibiting lower levels of phenotypic diversity. When the apical phase was being replaced by milli-Q, this led to a small dilution of the accumulated metabolites which stimulated the bacteria to become somewhat

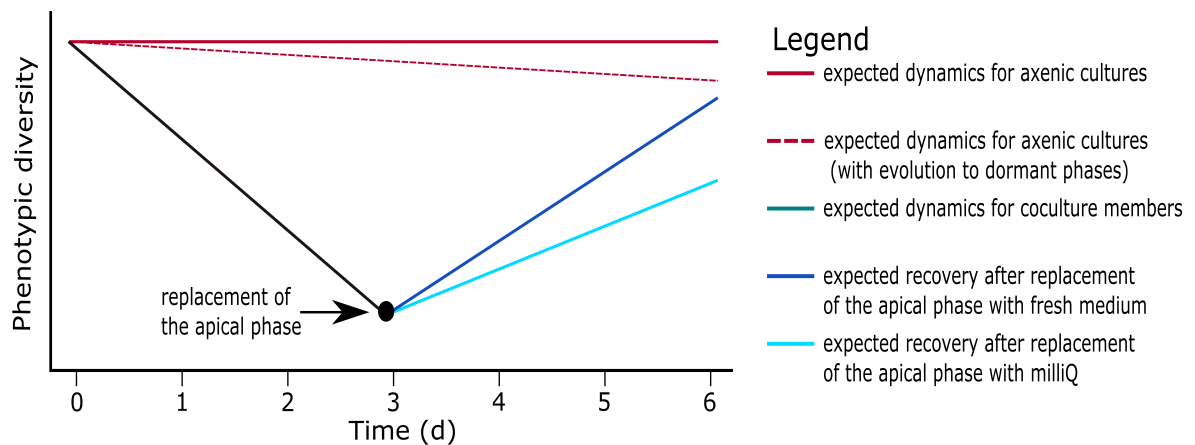


Figure 4.2: Hypothesis of the experiment where reversibility of the effect of interactions on phenotypic diversity was assessed. After three days of incubation, the apical phases of the cocultures were replaced with new apical phases, containing either fresh medium or milli-Q. The cultures where fresh medium was added were expected to grow and thus return to a higher level of phenotypic diversity at a fast rate (dark blue line). The replacement with milli-Q was intended to simulate the termination of any interaction between the genotypes. A recovery to higher levels of phenotypic diversity was expected after removing the interacting partner, but at a slower rate compared to the cultures where fresh medium was added (light blue line). The axenic cultures were intended as control samples. For these cultures either relatively stable levels of diversity were expected (red line) or a small decrease since the cultures were becoming less active (dashed red line).

more active again and therefore display a higher level of phenotypic diversity. We hypothesize that same could have happened when replacing the apical phase by fresh medium. For species B, the differences between the treatments were more limited, with a small increase in diversity for the cultures that were treated with milli-Q. These observations demonstrate that species-specific phenotypic responses can occur as a result of changing environments.

Influence of carbon source diversity

In the last experiment we assessed whether carbon source diversity drives phenotypic diversification in synthetic communities. An experiment was conducted where either glucose, a mixture of glucose, acetate and pyruvate or yeast extract was added to minimal medium. Our hypothesis stated that more carbon sources would lead to more niches in the community and therefore a higher level of phenotypic diversity would be observed. The diversity dynamics that were observed were related to community composition rather than to the carbon source that was added to the community (Figure 3.15 and Figure 3.16). This can potentially be explained by the fact that

all species were assumed to be in stationary phase and were therefore under starvation conditions. Bacterial strains are found to maintain a recognisable strain-specific DNA pattern when grown under limited conditions^[113]. If the cultures were also displaying similar RNA patterns because they were under starvation conditions, both DNA and RNA patterns would be similar for a certain genotype, irrespective of the carbon source that was added. This would explain why the axenic cultures were showing similar phenotypic diversity patterns under all conditions. Based on trends observed through beta-diversity analysis, the community composition was evolving in a similar way irrespective of the carbon sources (graph not shown). If both community composition and the influence of interaction on phenotypic diversity of the community members were similar, the community fingerprint would also evolve in a similar way. And as such result in similar diversity patterns for the mixed cultures. Potentially, the expected differences in diversity would have been observable throughout the growth of the bacteria, when they would not be under starvation. This could be evaluated by following up the phenotypic diversity during growth. Or, since bacteria exhibit different phenotypes during the different growth phases and the evolution of the growth phases might be different when different carbon sources are available, it might be difficult to distinguish between phenotypic diversity differences which are due to growth and differences which are due to the presence of multiple carbon sources. Therefore a chemostat set-up where growth rates can be controlled might be a more appropriate set-up to evaluate the effect of these types of influences on phenotypic diversity.

4.2 Hypothesis 2: Flow cytometry and Raman spectroscopy give complementary information regarding phenotypic community structure.

Phenotypic diversity is a population property that is manifesting itself at the level of individuals. When we want to assess this fine scale diversity we need tools that can reliably measure characteristics of single cells without disturbing their biochemical state^[114]. During this study phenotypes were evaluated using two techniques, flow cytometry and Raman spectroscopy. One of the aims in this study was to evaluate whether these techniques give similar or complementary information.

When we want to assess phenotypic diversity we need to define the phenotypes between which we will distinguish. For evaluating diversity based on flow cytometrically derived phenotypic traits, we use our in-house pipeline^[63] (explained in Section 1.4.2.1), where a binning grid is applied to each of the bivariate parameter combinations. Bacteria that fall within the same bin are defined as the same phenotype. Thus, phenotypes are defined in an arbitrary way, which implies that also our calculated diversity metrics, especially the phenotypic richness (D_0) which is the amount of non-empty bins, have somewhat arbitrary values. This arbitrary definition of phenotypes can be disputed. However, when diversity at the species level is assessed through sequencing, arbitrary cut-offs are used as well. The reads that result from sequencing are clustered together in OTU's (operational taxonomic units), which are used as a proxy for species. To define the OTU's an arbitrary cut-off of 97% similarity is frequently used, which is something that has been disputed in literature as well^[115]. However, we are not interested in the absolute values of the diversity metrics. It is rather by comparison of diversity values under different conditions or by following up the diversity over time, that we can gain insights in the underlying ecological processes that are occurring^[116].

When evaluating phenotypic diversity based on flow cytometry the phenotypic traits on which information is gained are cell size, morphology, granularity and nucleic acid content. Information regarding these traits is embedded in the scatter and fluorescence parameters. But, only a certain level of information is retained in these parameter representations^[117]. For example, based on the values of the scatter parameters the exact cell morphology (e.g. bacillus, rod, spiral, etc.) cannot directly be inferred. Moreover, since SG is a stain which stains all nucleic acids, it is unknown whether the fluorescence signal arises from DNA or RNA. Thus the phenotypic traits derived through flow cytometry give an abstract representation of the phenotypic traits. Moreover, only taking into account these traits is an abstraction of the entire phenotypic diversity of the

bacteria. Nonetheless, morphology and nucleic acids content are phenotypic traits of interest. DNA gives information regarding presence of species (e.g. genome sizes) and copy numbers of their genome. Previously flow cytometry-derived phenotypic diversity has been found to serve as a good proxy for taxonomic diversity in freshwater communities^[63]. Note that in this study the presence of species was a controlled factor. Additionally, the RNA content inferred from FCM gives information on activity of the bacteria. DNA and RNA levels are thus both properties of interest when assessing community dynamics regarding functionality. Even though the parameter representations are abstractions of the phenotypic traits, they can serve as comparative measures that are implicitly holding information regarding the phenotypic traits.

Heterogeneity is encountered for more cell constituents than nucleic acids alone. There can be both quantitative and qualitative differences in molecular phenotype between individuals, i.e. differences in concentration of biomolecules and differences in types of biomolecules^[34]. The Raman spectrum of a single cell is comprised of the spectra of all compounds that are present in the cell. Moreover, the signal intensity is proportional to the concentration of the compounds^[66], and therefore the spectrum offers an in depth view on the molecular phenotype. In Section 3.1.5 we attempted to gain some insight in whether specific biomolecules could be associated with the shifts in phenotypic diversity that were observed through flow cytometry. The inference of biochemical composition of bacterial cells based on their Raman spectra is not an automated procedure, and is therefore a time-consuming and unprecise task. In literature several approaches are used to gain insight in the spectra; some studies use visual comparison^[77] while others prefer the use of peak detection algorithms and subsequent visual comparison^[76]. A study of Bergholt *et al.*^[118] attempted to automate biomolecule detection from Raman spectra in a clinical set-up, however they focused on changes in only a few cancer-specific biomolecules and not on whole cell analysis. More automated pipelines for detection of biomolecules would be an interesting improvement.

Even though interpretation of the spectra is not straightforward, we argue the spectra can still be used for phenotypic diversity estimation. As is the case for flow cytometric parameter representations, an exact interpretation of the underlying phenotypic differences is not necessary to be able to use the spectra for comparative diversity analysis. We are not aware of any attempts to characterize phenotypic community structure based on Raman spectroscopy in literature. A first step to assess phenotypic diversity based on Raman spectra will be to find a way to define the phenotypes between which will be distinguished. During this study 51 single cell spectra were acquired for each population, without biological replicates. As indicated by the standard deviations that were found when calculating the average spectra for each population, there is a broad range of diversity

between cells of the same population (Figure 3.9). In a study of Nichols *et al.*^[119] single knock-out strains of *Escherichia coli* were subjected to a broad range of stresses in order to define the number of phenotypes that could be exhibited by each strain. The researchers concluded that the expected number of phenotypes ranged between 1 and 31. Therefore, we believe the 51 cells measured during this study do not provide enough sampling depth (i.e. enough coverage of the phenotypic landscape) for reliable diversity estimation.

Single-cell Raman measurements are time-consuming since they require manual focusing of the cells. The required exposure times are relatively long since Raman scatter is a weak signal^[120]. This requires an immobilization of the cells. The time consuming measurements cause the need for sample fixation and make it difficult to analyse large cell numbers or to replicate. In the field of Raman microscopy several methods have been developed to make the measurements shorter and more automated. For example through the use of microfluidics systems or by measuring single cell spectra directly in aqueous solution through optical trapping^[121]. Application of such techniques might provide a way to increase sampling depth and reduce the potentially human induced differences in the spectra through manual focusing etc.

In summary, flow cytometry is a high-throughput method for which we have an established diversity estimation pipeline. The main benefits of the flow cytometric approach are its speed and the fact that large amounts of cells can be analysed. This allows to have good coverage of the phenotypic landscape of the community and to apply a high measurement frequency. The properties which are taken into account by flow cytometric diversity are limited and thus the obtained diversity estimates are a coarse approximation of the entire diversity of the community. When a shift in phenotypic diversity is observed, it is not straightforward to draw a conclusion about the underlying biological or ecological process that is occurring. Raman spectroscopy on the other hand is a more sensitive technique that allows for a more holistic view on molecular phenotypic traits. This technique holds a lot of potential for phenotypic diversity estimation. However, there are some potential enhancements which would make the application easier, such as more automated and faster measurements since speed is currently the major bottleneck. The spectra are holding a lot of information, but they are difficult to interpret. More automated pipelines for detection of biomolecules would be an interesting improvement. Finally, it would be interesting to find confirmation of the phenotypic diversity estimates by flow cytometry through Raman spectroscopy. Considering the current study, both techniques have their benefits and drawbacks. They can give information on different (e.g. morphology, protein content) and similar (e.g. nucleic acids) properties, and therefore a combination of both technologies is interesting for further research.

4.3 Application of the in silico methodology to infer community composition

The supervised in silico community methodology described in Section 1.5, makes use of the cytometric fingerprints of species to infer community composition in synthetic ecology experiments. One of the aims in this study was to apply this newly developed method to a synthetic ecosystem study and to evaluate its performance.

In the first experiment, a shift in the cytometric fingerprint over time was observed for both species. When the species were grown in a coculture, where they could interact, their cytometric fingerprints were different from those of the axenic cultures. To evaluate whether these shifts, both in time and due to the interaction, in fingerprint would affect the model predictions, classifiers were trained in three ways using different input data (Figure 3.6). Based on the PCoA ordination the mixed culture was shifting from a culture that was resembling more to species A at 24h to a culture that was resembling more to species B at 48h and 72h; therefore, a gradual enrichment of species B in the mixed culture was expected (Figure 3.4). To evaluate the effect of the trend in time, predictions were made using a classifier that was trained on data of the axenic cultures from only the first timepoint or data of the axenic cultures for each of the corresponding timepoints. The predictions for these two approaches were very similar (differences $<1\%$). This might be explained by the fact that both species were in stationary phase, and thus even though there were some dynamics in their fingerprints the decision boundary did not change very much. To evaluate the effect of the changing fingerprints due to interaction, the relative abundances in the mixed culture were predicted by a classifier that was trained on the data of the membrane separated cultures as well. This resulted in very different predicted relative abundances compared to when the model was trained on data of the axenic cultures. At 24h, the differences in the predicted relative abundances were rather limited ($<10\%$). At 48h and 72h the differences were larger. The predictions based on the axenic cultures predicted that species A was present in a relative abundance of less than 1%, while the predictions based on the cocultures predicted a relative abundance of species A of about 30%. The fact that the differences in predicted relative abundance were more limited at 24h, corresponds with the previous observation that the differences in phenotypic fingerprint between the coculture and the axenic culture were limited at 24h and became larger over time (Figure 3.3). As the coculture data best represents the phenotypic behaviour of both species during interaction, we argue that the model trained by this dataset is the most biologically accurate.

To validate whether this conclusion was correct, an independent experiment was carried out where mixed cultures were created using a *gfp*-labeled strain of species A. This causes species A to emit green fluorescence and thus to be detectable on the first fluorescence channel of the flow cytometer without staining. The mixed cultures were analysed without staining to get information on the cell density of species A, and with SG staining to get information on the total cell density. This way the relative abundances in the mixed culture could be inferred. It should be noted that the presence of the *gfp*-label causes a metabolic burden for bacteria, which implies the relative abundances during the first experiment might be slightly diverging from the ones found in the third experiment. Next to this, the *gfp*-expression in a clonal population is heterogeneous^[33;101], therefore some members of species A will have remained undetected causing an underestimation of the relative abundances of species A. With these two factors in mind, the relative abundances of species A were expected to be slightly lower than the ones that were predicted in the first experiment. In the third experiment we found a higher abundance of species A at 24h and a gradual enrichment of species B from 48h on (Figure 3.14). This trend is similar to what was found when training the classifier on data of the cocultures. However, the exact values of the relative abundances were different. This could be due to differences in the experimental set-up. In the third experiment, the cultures were grown in 10mL tubes instead of 6-well plates and were not shaken while they were shaken during the first experiment. Competition is known to be influenced by mixing of cultures in laboratory settings^[122]. Therefore we believe this observed discrepancy of relative abundances in the first and third experiment might be explained by the differences in the experimental set-up. Since the trend in relative abundances in both experiments is similar and is confirmed by what was expected based on the beta-diversity analysis of the cultures (Figure 3.4), we believe the *in silico* communities are a reliable tool for inferring community compositions in a synthetic ecology experiment, provided that the correct (i.e. of interacting genotypes) input data is used to train the classifier. To have a final confirmation of the method, the coculture experiment could be repeated using a *gfp*-labeled strain.

Further, we applied the *in silico* methodology to validate our experimental set-up (Section 3.1.3). We created communities digitally, by merging and splitting data of measured cultures. Mixed cultures were created by merging data of species A and B. Coculture members were created by splitting data of the mixed cultures. By evaluating the similarity between the ordination based on the measurements and the ordination based on the predicted (*in silico*) populations, we aimed to validate whether the fingerprints of species A and species B in the coculture were the same as the fingerprints of species A and species B in the mixed culture. In other words, we wanted to validate

whether studying the effect of binary interactions on the phenotypic fingerprints in mixed cultures via a coculture set-up was a valid approach. In silico communities were created based on data of both the axenic cultures and the cocultures. Similarity of the ordinations was quantified using Procrustes analysis. Both in silico ordinations were significantly similar to the one of the measured samples (Table 3.1). Surprisingly, the ordination of the in silico created communities based on the axenic cultures was most similar to the one of the measured sample, however the difference was very small. A possible explanation for this might be that the evolution through time has a clear effect in the ordination, and thus even though there are some discrepancies in the interrelationships of the fingerprints the time effect dominates the sample positioning. Moreover, the lowest cell number in a population for the in silico created communities based on the axenic cultures was only 99 cells. To make a fair comparison all populations, including those of the measured samples, were subsampled to 99 cells. This probably caused an unreliable estimation of the phenotypic community structure with high estimated abundance for only the most abundant phenotypes. When these communities are then compared there is a comparison of the most abundant phenotypes and not of the entire community structure. Therefore we argue this comparison is probably unreliable. Despite that there are some small discrepancies when we visually compare the ordination of the measured samples and the ordination of the predicted (in silico) communities based on the cocultures, most interrelationships between the cultures are preserved well and the similarity between the cultures is significant. From this we conclude that studying the effect of binary interactions on the phenotypic fingerprints in mixed cultures via a coculture set-up was a valid approach.

4.4 Experimental set-up

There is a growing interest in understanding how phenotypic diversity is manifesting itself and what its potential importance might be in both natural and engineered microbial ecosystems. This interest arises from the growing awareness that bacterial heterogeneity is an essential trait for many biological processes^[101], such as pathogenicity^[123] and steering of microbial-based processes^[124]. It is still unknown how phenotypic diversity is influencing microbial communities and what the importance and implications of phenotypic diversity in natural environments might be^[33].

In literature, phenotypic diversity is often studied using phenotypic arrays^[31], isotope labeling with stable or radioactive probes or fluorescent labeled proteins^[101;103]. Via phenotypic arrays insights can be gained concerning the potential phenotypes that a certain species or an entire community can exhibit. These set-ups are used to evaluate diversity in for example metabolic potential, but cannot be used to quantify diversity within clonal populations in the same environment. Both isotope labeling and fluorescent labeled proteins do allow to study diversity in clonal populations. However, they require either a modification of the organisms under study by inserting a fluorescent protein or the use of rather expensive and potentially dangerous isotopes.

The experimental set-up applied in this study does not require any genetic alteration of bacteria or the use of isotopes. It can be used to study both biotic (interacting partner or interacting community) and abiotic (growing medium, temperature, agitation, perturbations, etc.) factors that might influence phenotypic diversity of microbes. It should be noted that this set-up does not account for interactions which involve physical contact. Nonetheless is it an interesting experimental set-up to study freshly isolated bacteria. Additional insights could be gained in the ecology of interesting bacteria, such as the genus *Limnohabitans*, discussed in Section 1.1.2. Potentially, their omnipresence might be related to their phenotypic diversity or phenotypic plasticity under different circumstances, such as competition.

Thanks to recent advances in the field of flow cytometry, high-frequency automated sampling is now possible as well^[125]. After some adaptations in the experimental set-up, for example the use of dialysis membranes to separate bacterial populations in larger volume cultures, it would be possible to extend this set-up to follow up microbial populations at a higher measurement frequency through online FCM. This could be interesting in the context of follow-up of growth dynamics in mixed cultures or to evaluate the response of microbial populations on perturbations. Experimental set-ups that allow these kinds of studies are of interest and are being developed in recent literature^[126].

4.5 Conclusion and further perspectives

The main finding of this study is that interactions between sympatric bacterial populations can lead to an adjustment of the individual phenotypic diversities of the interacting populations. This is an interesting finding which opens a future for research on phenotypic diversity and microbial interactions. This can be within the context of both engineered ecosystems, such as a refined steering and optimisation of microbial-based processes, or natural ecosystems, such as an increased understanding of the drivers and sensitivities of microbial communities in natural ecosystems. We proposed a hypothesis which stated that the individual phenotypic diversity of a genotype is function of the other genotypes in the community and that this might be related to niche separation between the genotypes. Our work has now provided the framework under which this hypothesis can be further evaluated at higher genotypic richness.

Furthermore, we concluded that the experimental set-up that was used in this study was very suitable for its purpose. It can be used to study both biotic and abiotic factors that might influence phenotypic diversity of microbes. We propose to extend this set-up for applications where a higher sampling frequency is desirable. This would allow to monitor microbial population dynamics at a very fine temporal resolution which will enhance our understanding of microbial community dynamics. Interesting experimental set-ups could be used to monitor the growth of single genotypes in mixed communities, which is currently a difficulty, or to follow-up microbial responses to perturbations.

We evaluated two techniques, flow cytometry, which had previously been used for phenotypic diversity estimation, and Raman spectroscopy, which has never been used for this purpose. We concluded that both techniques have their benefits and drawbacks, that they give similar and complementary information and that a combination is necessary. In this study we were not able to estimate phenotypic diversity based on Raman spectroscopy. Nonetheless, we were able to gain valuable information regarding its bottlenecks and its potential for evaluating phenotypic diversity. Furthermore, we suggest that a pipeline for a more automated detection of biomolecules would be an interesting enhancement for microbial research using Raman spectroscopy. This would complement the fast diversity screening of FCM.

Following up the community composition in synthetic ecosystem experiments is currently a difficulty. The available tools either require large sample volumes (sequencing) or the development of specific primers (qPCR). Through phenotypic fingerprinting we found that we are able to infer the evolution of the community composition using beta-diversity analysis, but more importantly, that

we can quantify the relative abundances in a mixed culture experimental set-up by means of supervised machine learning techniques. We concluded that *in silico* communities are a reliable tool for inferring community composition in a co-culture synthetic ecology experiment, provided that the correct input data is used to train the classifier. We advise to validate the approach by means of a *gfp*-labeled strain before applying the approach to new experiments. We also found that the phenotypic community structure of a species is a very dynamic property, in time as well as in relation to external influencing factors. And, that this causes the *in silico* approach to break-down in case of experimental set-ups which do not involve a co-culture standard. Therefore, research regarding development of the tools necessary for characterizing the community composition in synthetic ecology experiments is still ongoing.

In summary, in this thesis we found that interactions between sympatric bacterial populations can lead to an adjustment of the individual phenotypic diversities of the interacting populations. We evaluated the potential of Raman spectroscopy to estimate phenotypic diversity. The experimental design presented here forms a framework within which new ecological hypotheses regarding phenotypic diversity and microbial interactions can be tested.

Bibliography

- [1] Sigee D. *Freshwater Microbiology: Biodiversity and Dynamic Interactions of Microorganisms in the Aquatic Environment*. Wiley, 2005.
- [2] Bell T., Newman J.A., Silverman B.W., Turner S.L., and Lilley A.K. The contribution of species richness and composition to bacterial services. *Nature*, 436:1157–1160, 2005.
- [3] Wittebolle L., Marzorati M., Clement L., Balloi A., Daffonchio D., De Vos P., Heylen K., Verstraete W., and Boon N. Initial community evenness favours functionality under selective stress. *Nature*, 458:623–626, 2009.
- [4] Dolinšek J., Goldschmidt F., and Johnson D.R. Synthetic microbial ecology and the dynamic interplay between microbial genotypes. *FEMS Microbiology Reviews*, 112:961–979, 2016.
- [5] Klitgord N. and Segre D. Environments that Induce Synthetic Microbial Ecosystems. *PLoS Computational Biology*, 6:1–17, 2010.
- [6] Escalante A.E., Rebolleda-Gómez M., Benítez M., and Travisano M. Ecological perspectives on synthetic biology: Insights from microbial population biology. *Frontiers in Microbiology*, 6:1–10, 2015.
- [7] Travisano M. and Velicer G.J. Strategies of microbial cheater control. *Trends in Microbiology*, 12:72–78, 2004.
- [8] Begon M., Townsend C.R., and Harper J.L. *Ecology - From Individuals to Ecosystems*. Blackwell publishing, 2007.
- [9] Little A.E., Robinson C.J., Peterson S.B., Raffa K.F., and Handelsman J. Rules of Engagement: Interspecies Interactions that Regulate Microbial Communities. *Annu Rev Microbiol*, 62:375–401, 2008.
- [10] Porter J., Deere D., Pickup R., and Edwards C. Fluorescent probes and flow cytometry: New insights into environmental bacteriology. *Cytometry*, 23:91–96, 1996.
- [11] Newton R.J., Jones S.E., Eiler A., McMahon K.D., and Bertilsson S. A guide to the natural history of freshwater lake bacteria. *Microbiology and molecular biology reviews*, 75:14–49, 2011.

- [12] Pernthaler J. and Amann R.I. Fate of heterotrophic microbes in pelagic habitats: focus on populations. *Microbiology and Molecular Biology Reviews*, 69:440–461, 2005.
- [13] Lindström E. Bacterioplankton Community Composition in Five Lakes Differing in Trophic Status and Humic Content. *Microbial ecology*, 40:104–113, 2000.
- [14] Hahn M.W., Kasalický V., Jezbera J., Brandt U., Jezberová J., and Šimek K. *Limnohabitans curvus* gen. nov., sp. nov., a planktonic bacterium isolated from a freshwater lake. *International Journal of Systematic and Evolutionary Microbiology*, 60:1358–1365, 2010.
- [15] Kasalický V., Jezbera J., Hahn M.W., and Šimek K. The Diversity of the *Limnohabitans* Genus, an Important Group of Freshwater Bacterioplankton, by Characterization of 35 Isolated Strains. *PLoS ONE*, 8:1–13, 2013.
- [16] Andersson A., Larsson U., and Hagström Å. Size-selective grazing by a microflagellate on pelagic bacteria. *Marine Ecology - Progress Series*, 33:51–57, 1986.
- [17] Jezbera J., Jezberová J., Koll U., Horňák K., Šimek K., and Hahn M.W. Contrasting trends in distribution of four major planktonic betaproteobacterial groups along a pH gradient of epilimnia of 72 freshwater habitats. *FEMS Microbiology Ecology*, 81(2):467–479, 2012.
- [18] Jezbera J., Jezberová J., Kasalický V., Šimek K., and Hahn M.W. Patterns of *Limnohabitans* Microdiversity across a Large Set of Freshwater Habitats as Revealed by Reverse Line Blot Hybridization. *PLoS ONE*, 8:1–10, 2013.
- [19] Šimek K., Kasalický V., Jezbera J., Jezberová J., Hejzlar J., and Hahn M.W. Broad habitat range of the phylogenetically narrow R-BT065 cluster, representing a core group of the betaproteobacterial genus *limnohabitans*. *Applied and Environmental Microbiology*, 76:631–639, 2010.
- [20] De Roy K., Marzorati M., Van den Abbeele P., Van de Wiele T., and Boon N. Synthetic microbial ecosystems: an exciting tool to understand and apply microbial communities. *Environmental microbiology*, 16:1472–1481, 2014.
- [21] Großkopf T. and Soyer O.S. Synthetic microbial communities. *Current Opinion in Microbiology*, 18:72–77, 2014.
- [22] Jessup C.M., Kassen R., Forde S.E., Kerr B., Buckling A., Rainey P.B., and Bohannan B.J.M. Big questions, small worlds: Microbial model systems in ecology. *Trends in Ecology and Evolution*, 19:189–197, 2004.

- [23] Yu Z., Krause S.M.B., Beck D.A.C., and Chistoserdova L. A synthetic ecology perspective: How well does behavior of model organisms in the laboratory predict microbial activities in natural habitats? *Frontiers in Microbiology*, 7:1–7, 2016.
- [24] Goers L., Freemont P., and Polizzi K.M. Co-culture systems and technologies: taking synthetic biology to the next level. *Journal of the Royal Society*, 11:1–13, 2014.
- [25] Zengler K., Toledo G., Rappe M., Elkins J., Mathur E.J., Short J.M., and Keller M. Cultivating the uncultured. *Proceedings of the National Academy of Sciences of the United States of America*, 99:15681–15686, 2002.
- [26] De Ryck T., Grootaert C., Jaspert L., Kerckhof F.M., Van Gele M., De Schrijver J., Van Den Abbeele P., Swift S., Bracke M., Van De Wiele T., and Vanhoecke B. Development of an oral mucosa model to study host-microbiome interactions during wound healing. *Applied Microbiology and Biotechnology*, 98:6831–6846, 2014.
- [27] Dunham M.J. Synthetic ecology: a model system for cooperation. *Proceedings of the National Academy of Sciences of the United States of America*, 104:1741–1742, 2007.
- [28] Köhler J. Measurement of in situ growth rates of phytoplankton under conditions of simulated turbulence. *Journal of Plankton Research*, 19:849–862, 1997.
- [29] Schlegel H.G., Zaborosch C., and Kogut M. *General Microbiology*. Cambridge University Press, 1993.
- [30] Bochner B.R. Global phenotypic characterization of bacteria. *FEMS Microbiology Reviews*, 33:191–205, 2009.
- [31] Borglin S., Joyner D., DeAngelis K.M., Khudyakov J., D’haeseleer P., Joachimiak M.P., and Hazen T. Application of phenotypic microarrays to environmental microbiology. *Current Opinion in Biotechnology*, 23:41–48, 2012.
- [32] Turcotte M.M. and Levine J.M. Phenotypic Plasticity and Species Coexistence. *Trends in Ecology & Evolution*, pages 803–813, 2016.
- [33] Ackermann M. Microbial individuality in the natural environment. *The ISME Journal*, 7:465–467, 2013.
- [34] Ackermann M. A functional perspective on phenotypic heterogeneity in microorganisms. *Nature Publishing Group*, 13:497–508, 2015.

- [35] Schreiber F., Littmann S., Lavik G., Escrig S., Meibom A., Kuypers M.M.M., and Ackermann M. Phenotypic heterogeneity driven by nutrient limitation promotes growth in fluctuating environments. *Nature Microbiology*, in press:1–7, 2016.
- [36] Ackermann M. and Schreiber F. A growing focus on bacterial individuality. *Environmental Microbiology*, 17:2193–2195, 2015.
- [37] Taylor J. *Microorganisms and Biotechnology*. Nelson Thornes, 2001.
- [38] Kratz R. *Microbiology the Easy Way*. Barron's, 2005.
- [39] Pommerville J.C. *Alcama's Laboratory Fundamentals of Microbiology*. Jones and Bartlett, 2007.
- [40] Egli T. Microbial growth and physiology: A call for better craftsmanship. *Frontiers in Microbiology*, 6:1–12, 2015.
- [41] Rao D.G. *Introduction to Biochemical Engineering*. Tata McGraw-Hill, 2010.
- [42] Pepper I.L. and Gerba C.P. *Environmental Microbiology: A Laboratory Manual*. Elsevier Academic Press, 2005.
- [43] Monod J. The Growth of Bacterial Cultures. *Annual Review of Microbiology*, 3:371–394, 1949.
- [44] Wang L., Fan D., Chen W., and Terentjev E.M. Bacterial growth, detachment and cell size control on polyethylene terephthalate surfaces. *Scientific Reports*, 5:1–11, 2015.
- [45] Cooper S. *Bacterial Growth and Division: Biochemistry and Regulation of Prokaryotic and Eukaryotic Division Cycles*. Elsevier Science, 2012.
- [46] Koch A. *Bacterial Growth and Form*. Springer Netherlands, 2013.
- [47] Cooper S. Bacterial Growth and Division. *Encyclopedia of Molecular Cell Biology and Molecular Medicine*, pages 1–27, 2006.
- [48] BDBiosciences. *Introduction to Flow Cytometry: A Learning Guide*. 2000.
- [49] De Roy K. *Microbial Resource Management: Introducing New Tools and Ecological Theories*. PhD thesis, Ghent University, Belgium, 2014.

- [50] Müller S. and Nebe-Von-Caron G. Functional single-cell analyses: Flow cytometry and cell sorting of microbial populations and communities. *FEMS Microbiology Reviews*, 34:554–587, 2010.
- [51] Hammes F. and Egli T. Cytometric methods for measuring bacteria in water: Advantages, pitfalls and applications. *Analytical and Bioanalytical Chemistry*, 397:1083–1095, 2010.
- [52] Gatza E., Hammes F., and Prest E. *Assessing Water Quality with the BD Accuri C6 Flow Cytometer*. BDBiosciences, 2009.
- [53] Nebe-Von-Caron G., Stephens P.J., Hewitt C. J., Powell J.R., and R. A. Badley. Analysis of bacterial function by multi-colour fluorescence flow cytometry and single cell sorting. *Journal of Microbiological Methods*, 42:97–114, 2000.
- [54] Ormerod M.G. *Flow Cytometry: A Practical Approach*. OUP Oxford, 2000.
- [55] Hyka P., Lickova S., Přibyl P., Melzoch K., and Kovar K. Flow cytometry for the development of biotechnological processes with microalgae. *Biotechnology Advances*, 31:2–16, 2013.
- [56] Barbesti S., Citterio S., Labra M., Baroni M.D., Neri M.G., and Sgorbati S. Two and three-color fluorescence flow cytometric analysis of immunoidentified viable bacteria. *Cytometry*, 40:214–218, 2000.
- [57] Grégori G., Citterio S., Ghiani A., Labra M., Sgorbati S., Brown S, and Denis M. Resolution of Viable and Membrane-Compromised Bacteria in Freshwater and Marine Waters Based on Analytical Flow Cytometry and Nucleic Acid Double Staining. *Applied and Environmental Microbiology*, 67:4662–4670, 2001.
- [58] Rubbens P., Props R., Boon N., and Waegeman W. Flow cytometric single-cell identification of populations in synthetic bacterial communities. *PLOS ONE*, 12:1–19, 2017.
- [59] Van Nevel S., Koetzsch S., Weilenmann H.U., Boon N., and Hammes F. Routine bacterial analysis with automated flow cytometry. *Journal of Microbiological Methods*, 94:73–76, 2013.
- [60] Epstein S.S. The phenomenon of microbial uncultivability. *Current Opinion in Microbiology*, 16:636–642, 2013.
- [61] Shapiro H.M. *Practical Flow Cytometry*. Wiley, 2005.

- [62] Bashashati A. and Brinkman R.R. A Survey of Flow Cytometry Data Analysis Methods. *Advances in Bioinformatics*, 2009:1–19, 2009.
- [63] Props R., Monsieurs P., Mysara M., Clement L., and Boon N. Measuring the biodiversity of microbial communities by flow cytometry. *Methods in Ecology and Evolution*, 7:1376–1385, 2016.
- [64] Hastings A. and Gross L. *Encyclopedia of Theoretical Ecology*. University of California Press, 2012.
- [65] Greenacre M. and Primicerio R. *Multivariate Analysis of Ecological Data*. Fundación BBVA, 2014.
- [66] Larkin P. *Infrared and Raman Spectroscopy: Principles and Spectral Interpretation*. Elsevier Science, 2011.
- [67] Koenig J.L. *Infrared and Raman Spectroscopy of Polymers*. Rapra Technology, 2001.
- [68] B&W TEK. Theory of Raman Scattering, 2016.
- [69] Ferraro J.R. and Nakamoto K. *Introductory Raman Spectroscopy*. Elsevier Science, 2012.
- [70] Dieing T., Hollricher O., and Toporski J. *Confocal Raman Microscopy*. Springer Berlin Heidelberg, 2011.
- [71] Huang W.E., Li M., Jarvis R.M., Goodacre R., and Banwart S.A. *Shining light on the microbial world the application of Raman microspectroscopy*. Elsevier Inc., 2010.
- [72] Kaiser Optical Systems. Raman Tutorial, 2016.
- [73] Butler H.J., Ashton L., Bird B., Cinque G., Curtis K., Esmonde-white K., Fullwood N.J., Gardner B., Martin-hirsch P.L., Walsh M.J., McAinsh M.R., Stone N., and Martin F.L. Using Raman spectroscopy to characterise biological materials. *Nature Protocols*, 11:1–47, 2016.
- [74] Schrader B. *Infrared and Raman Spectroscopy: Methods and Applications*. Wiley, 2008.
- [75] van de Vossenberg J., Tervahauta H., Maquelin K., Blokker- Koopmans C.H.W., Uytewaal-Aarts M., van der Kooij D., VanWezel A.P., and van der Gaag B. Identification of bacteria in drinking water with Raman. *Analytical Methods*, 5:2679–2687, 2013.

- [76] Read D.S., Woodcock D.J., Strachan N.J.C., Forbes K.J., Colles F.M., Maiden M.C.J., Clifton-hadley F., Ridley A., Vidal A., Rodgers J., Whiteley A.S., and Sheppard K. Evidence for Phenotypic Plasticity among Multihost *Campylobacter jejuni* and *C. coli* Lineages, Obtained Using Ribosomal Multilocus Sequence Typing and Raman Spectroscopy. *Applied and Environmental Microbiology*, 79:965–973, 2013.
- [77] De Gelder J., De Gussem K., Vandenabeele P., and Moens L. Reference database of Raman spectra of biological molecules. *Journal of Raman Spectroscopy*, 38:1133–1147, 2007.
- [78] Berry D., Mader E., Lee T.K., Wobken D., Wang Y., Zhu D., Palatinszky M., Schintlmeister A., Schmid M.C., Hanson B.T., Shterzer N., Mizrahi I., Rauch I., Decker T., Bocklitz T., Popp J., Gibson C.M., Fowler P.W., Huang W.E., and Wagner M. Tracking heavy water (D₂O) incorporation for identifying and sorting active microbial cells. *Proceedings of the National Academy of Sciences of the United States of America*, 112:194–203, 2015.
- [79] Stöckel S., Kirchhoff J., Neugebauer U., Rösch P., and Popp J. The application of Raman spectroscopy for the detection and identification of microorganisms. *Journal of Raman Spectroscopy*, 47:89–109, 2015.
- [80] Bocklitz T., Walter A., Hartmann K., Rösch P., and Popp J. How to pre-process Raman spectra for reliable and stable models? *Analytica Chimica Acta*, 704:47–56, 2011.
- [81] Rajwa B., Venkatapathi M., Ragheb K., Banada P.P., Hirleman E.D., Lary T., and Robinson J.P. Automated classification of bacterial particles in flow by multiangle scatter measurement and support vector machine classifier. *Cytometry Part A*, 73:369–379, 2008.
- [82] Davey H.M., Jones A., Shaw A.D., and Kell D.B. Variable selection and multivariate methods for the identification of microorganisms by flow cytometry. *Cytometry*, 35:162–168, 1999.
- [83] Boddy L., Morris C.W., Wilkins M.F., Tarran G.A., and Burkill P.H. Neural network analysis of flow cytometric data for 40 marine phytoplankton species. *Cytometry*, 15:283–293, 1994.
- [84] Pereira G.C. and Ebecken N.F.F. Combining in situ flow cytometry and artificial neural networks for aquatic systems monitoring. *Expert Systems with Applications*, 38:9626–9632, 2011.

- [85] Vilchez-Vargas R., Geffers R., Suárez-Diez M., Conte I., Waliczek A., Kaser V.S., Kralova M., Junca H., and Pieper D.H. Analysis of the microbial gene landscape and transcriptome for aromatic pollutants and alkane degradation using a novel internally calibrated microarray system. *Environmental Microbiology*, 15:1016–1039, 2013.
- [86] Nair A.J. *Principles of Biotechnology*. Laxmi Publications, 2008.
- [87] Zwietering M., Jongenburger I., Rombouts F., and Van't Riet K. Modeling of the Bacterial Growth Curve. *Applied and Environmental Microbiology*, 56:1875–1881, 1990.
- [88] Kniggendorf A., Gaul T., and Meinhardt-wollweber M. Effects of Ethanol, Formaldehyde, and Gentle Heat Fixation in Confocal Resonance Raman Microscopy of Purple NonSulfur Bacteria. *Microscopy Research and Technique*, 183:177–183, 2011.
- [89] R Core Team. *R: A Language and Environment for Statistical Computing*, 2016.
- [90] Kahm M., Hasenbrink G., Lichtenberg-Frat'e H., Ludwig J., and Kschischo M. profit: Fitting biological growth curves with R. *Journal of Statistical Software*, 33:1–21, 2010.
- [91] Ellis B., Haaland P., Hahne F., Le Meur N., Gopalakrishnan N., Spidlen J., and Jiang M. *flowCore: flowCore: Basic structures for flow cytometry data*, 2016.
- [92] Fletez-Brant K., Špidlen J., Brinkman R.R., Roederer M., and Chattopadhyay P.K. flow-Clean: Automated identification and removal of fluorescence anomalies in flow cytometry data. *Cytometry Part A*, 89:461–471, 2016.
- [93] 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 É. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2012.
- [94] Gibb S. and Strimmer K. MALDIquant: a versatile R package for the analysis of mass spectrometry data. *Bioinformatics*, 28:2270–2271, 2012.
- [95] Ryan C.G., Clayton E., Griffin W.L., Sie S.H., and Cousens D.R. SNIP, a statistics-sensitive background treatment for the quantitative analysis of PIXE spectra in geoscience applications. *Nuclear Inst. and Methods in Physics Research, B*, 34:396–402, 1988.
- [96] Rubbens P., Props R., Garcia C., Boon N., and Waegeman W. Stripping flow cytometry: how many detectors do we need for bacterial identification? (*under review*).

- [97] Baetens V. Monitoring of microbial communities in freshwater systems by flow cytometry. Master's thesis, Ghent University, Belgium, 2016.
- [98] Swinnen I.A.M., Bernaerts K., Dens E.J.J., Geeraerd A.H., and Van Impe J.F. Predictive modelling of the microbial lag phase: A review. *International Journal of Food Microbiology*, 94:137–159, 2004.
- [99] ThermoFisher Scientific. *The Molecular Probes Handbook: A Guide to Fluorescent Probes and Labeling Technologies*. Life Technologies Corporation, 2010.
- [100] Lieder S., Jahn M., Koepff J., Müller S., and Takors R. Environmental stress speeds up DNA replication in *Pseudomonas putida* in chemostat cultivations. *Biotechnology Journal*, 11:155–163, 2016.
- [101] Elowitz M.B., Levine A.J., and Siggia E.D. Stochastic Gene Expression in a Single Cell. *Science*, 297:1183–1186, 2002.
- [102] Ansel J., Bottin H., Rodriguez-Beltran C., Damon C., Nagarajan M., Fehrmann S., Francois J., and Yvert G. Cell-to-cell stochastic variation in gene expression is a complex genetic trait. *PLoS Genetics*, 4:1–10, 2008.
- [103] Newman J.R.S., Ghaemmaghami S., Ihmels J., Breslow D.K., Noble M., DeRisi J.L., and Weissman J.S. Single-cell proteomic analysis of *S. cerevisiae* reveals the architecture of biological noise. *Nature*, 441:840–846, 2006.
- [104] Fraser D. and Kærn M. A chance at survival: Gene expression noise and phenotypic diversification strategies. *Molecular Microbiology*, 71:1333–1340, 2009.
- [105] Trefil J.S. *Encyclopedia of Science and Technology*. Routledge, 2001.
- [106] Kuypers M.M.M. and Barker Jørgensen B. The future of single-cell environmental microbiology. *Environmental Microbiology*, 9:1–11, 2007.
- [107] Kassen R., Llewellyn M., and Rainey P.B. Ecological constraints on diversification in a model adaptive radiation. *Letters to Nature*, 431:984–988, 2004.
- [108] Grudemo J. and Johannesson K. Size of mudsnails, *Hydrobia ulvae* (Pennant) and *H. ventrosa* (Montagu), in allopatry and sympatry: Conclusions from field distributions and laboratory growth experiments. *Journal of Experimental Marine Biology and Ecology*, 239:167–181, 1999.

- [109] Anderson G.R.V., Ehrlich A.H., Ehrlich P.R., Roughgarden J.D., Russell B.C., and Talbot H. The Community Structure of Coral Reef Fishes. *The American Naturalist*, 117:476–495, 2008.
- [110] Smith V.H. Microbial diversity-productivity relationships in aquatic ecosystems. *FEMS Microbiology Ecology*, 62:181–186, 2007.
- [111] Pozo M.I., Herrera C.M., Lachance M.A., Verstrepen K., Lievens B., and Jacquemyn H. Species coexistence in simple microbial communities: unravelling the phenotypic landscape of co-occurring *Metschnikowia* species in floral nectar. *Environmental Microbiology*, 18:1850–1862, 2016.
- [112] Narwani A., Bentlage B., Alexandrou M.A., Fritschie K.J., Delwiche C., Oakley T.H., and Cardinale B.J. Ecological interactions and coexistence are predicted by gene expression similarity in freshwater green algae. *Journal of Ecology*, 105:580–591, 2017.
- [113] Müller S. and Babel W. Analysis of bacterial DNA patterns - An approach for controlling biotechnological processes. *Journal of Microbiological Methods*, 55:851–858, 2003.
- [114] Spudich J.L. and Koshland D.E. Non-genetic individuality: chance in the single cell. *Nature*, 262:467–471, 1976.
- [115] Mysara M., Vandamme P., Props R., Kerckhof F.M., Leys N., Boon N., Raes J., and Monsieus P. Reconciliation between operational taxonomic units and species boundaries. *FEMS Microbiology Ecology*, 93:1–12, 2017.
- [116] Shade A. Diversity is the question, not the answer. *The ISME Journal*, 4:1–6, 2017.
- [117] Fontana S., Jokela J., and Pomati F. Opportunities and challenges in deriving phytoplankton diversity measures from individual trait-based data obtained by scanning flow-cytometry. *Frontiers in Microbiology*, 5:1–12, 2014.
- [118] Bergholt M.S., Zheng W., Lin K., Ho K.Y., Teh M., Yeoh K.G., So J.B.Y., and Huang Z. In Vivo Diagnosis of Esophageal Cancer Using Image-Guided Raman Endoscopy and Biomolecular Modeling. *Technology in Cancer Research and Treatment*, 10:103–112, 2011.
- [119] Nichols R.J., Sen S., Choo Y.J., Beltrao P., Zietek M., Chaba R., Lee S., Kazmierczak K.M., Lee K.J., Wong A., Shales M., Lovett S., Winkler M.E., Krogan N.J., Typas A., and Gross C.A. Phenotypic landscape of a bacterial cell. *Cell*, 144:143–156, 2011.

- [120] Ando J., Palonpon A.F., Sodeoka M., and Fujita K. High-speed Raman imaging of cellular processes. *Current Opinion in Chemical Biology*, 33:16–24, 2016.
- [121] Xie C., Mace J., Dinno M.A., Li Y.Q., Tang W., Newton R.J., and Gemperline P.J. Identification of single bacterial cells in aqueous solution using confocal laser tweezers Raman spectroscopy. *Analytical Chemistry*, 77:4390–4397, 2005.
- [122] Hibbing M.E., Fuqua C., Parsek M.R., and Peterson S.B. Bacterial competition: surviving and thriving in the microbial jungle. *National Review of Microbiology*, 8:15–25, 2010.
- [123] Ackermann M., Stecher B., Freed N.E., Songhet P., Hardt W.D., and Doebeli M. Self-destructive cooperation mediated by phenotypic noise. *Nature*, 454:987–990, 2008.
- [124] Lencastre Fernandes R., Nierychlo M., Lundin L., Pedersen A.E., Puentes Tellez P.E., Dutta A., Carlquist M., Bolic A., Schäpper D., Brunetti A.C., Helmark S., Heins A.L., Jensen A.D., Nopens I., Rottwitt K., Szita N., van Elsas J.D., Nielsen P.H., Martinussen J., Sørensen S.J., Lantz A.E., and Gernaey K.V. Experimental methods and modeling techniques for description of cell population heterogeneity. *Biotechnology Advances*, 29:575–599, 2011.
- [125] Hammes F., Broger T., Weilenmann H.U., Vital M., Helbing J., Bosshart U., Huber P., Peter Odermatt R., and Sonnleitner B. Development and laboratory-scale testing of a fully automated online flow cytometer for drinking water analysis. *Cytometry Part A*, 81:508–516, 2012.
- [126] Moutinho T.J., Panagides J.C., Biggs M.B., Medlock G.L., Kolling G.L., and Papin J.A. Novel co-culture plate enables growth dynamic-based assessment of contact-independent microbial interactions. *doi:10.1101/145615*, 2017.
- [127] Anderson M.J. A new method for non parametric multivariate analysis of variance. *Austral ecology*, 26:32–46, 2001.
- [128] Jackson D.A. PROTEST: A PROcrustean Randomization TEST of community environment concordance. *Ecoscience*, 2:297–303, 1995.
- [129] Breiman L. Random Forests. *Machine Learning*, 45:1–33, 2001.
- [130] James G., Witten D., Hastie T., and Tibshirani R. *An Introduction to Statistical Learning*. Springer, 2013.

-
- [131] Meinshausen N. and Bühlmann P. Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72:417–473, 2010.
- [132] Jolliffe I.T. *Principal Component Analysis*. Springer New York, 2013.
- [133] Zuur A., Ieno E.N., and Smith G.M. *Analyzing Ecological Data*. Springer New York, 2007.
- [134] Digby P.G.N. and Kempton R.A. *Multivariate Analysis of Ecological Communities*. Springer Netherlands, 2012.
- [135] Van Der Maaten L.J.P. and Hinton G.E. Visualizing high-dimensional data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008.

5.1 Supplementary information data analysis

5.1.1 Bray-Curtis dissimilarity

The Bray-Curtis dissimilarity is widely used in ecology for comparison of species abundance data^[65] and is given by Equation 5.1.

$$D_{AB} = \frac{\sum_{i=1}^S |p_{Ai} - p_{Bi}|}{\sum_{i=1}^S (p_{Ai} + p_{Bi})}. \quad (5.1)$$

A simple illustration of the Bray-Curtis dissimilarity applied to the phenotypic fingerprints is given in Figure 5.1. Note that Bray-Curtis dissimilarity is a ‘dissimilarity’ and not a ‘distance’. This is because it does not satisfy the triangle inequality. The triangle inequality states that the sum of the lengths of any two sides of a triangle is greater than the length of the third side ($D_{AB} \leq D_{AC} + D_{CB}$), which is not valid in Bray-Curtis space^[65].

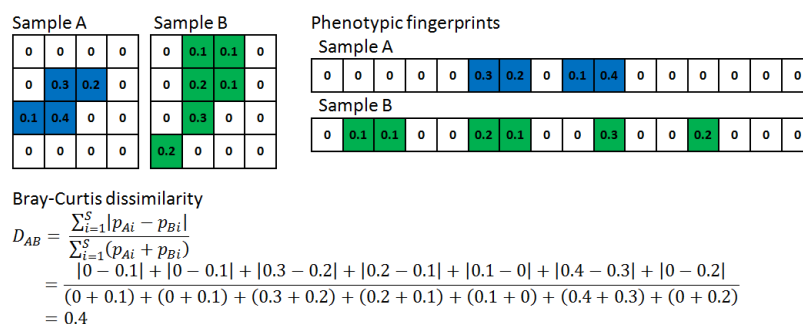


Figure 5.1: Simplified illustration of the Bray-curtis dissimilarity on phenotypic fingerprints. In the upper left corner a two dimensional fingerprint for two samples, A and B, after density estimation and normalisation. These are converted into phenotypic fingerprints by concatenating all bins into a one-dimensional vector. For these vectors the Bray-Curtis dissimilarity can be calculated using Equation 5.1.

5.1.2 Permutational multivariate analysis of variance using distance matrices (permanova)

Permanova^[127] is a non-parametric test to analyze variance for multivariate data. In a classical analysis of variance approach for univariate data the total sum of squares (SS_T) is partitioned in the within-group sum of squares (SS_W , i.e. sums of squared differences between individual replicates and their group mean) and the among-group sum of squares (SS_A , i.e. sums of squared differences between group means and the overall sample mean). From the ratio of SS_A over SS_W , we can see whether it is likely that the null hypothesis, which states there are no significant differences between the group means, is false. This concept can be expanded to the multivariate setting. SS_W can be thought of as the sum of the squared distances of each point to the centroid of the group to which the points belong, while SS_A is the sum of squared distances from group centroids to the overall centroid (Figure 5.2 A).

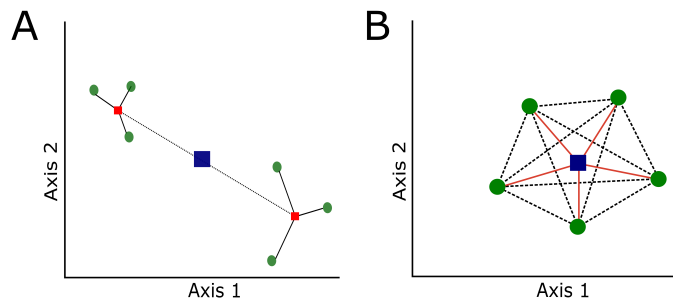


Figure 5.2: **A:** In a multivariate setting SS_W can be thought of as the sum of the squared distances of each point (green dots) to the centroid of the group (red squares) to which the points belong. SS_A is the sum of squared distances from group centroids (red squares) to the overall centroid (blue square). **B:** The summed distances from points to their group centroids (red lines) is equal to the summed inter-point distances (grey dotted lines) divided by the number of points in this group, 5 in this case. Figures redrafted after Anderson, 2001^[127].

Suppose we have a multivariate dataset with a groups, each containing n replicates. This dataset can be presented as a matrix where each row is a sample and each column is a variable (Figure 5.3 A). To evaluate whether the a groups are different, the values of SS_W and SS_A are needed. However, when similarity between these groups is evaluated based on some similarity or dissimilarity metric, finding the centroids might be problematic. For example, if the Bray-Curtis dissimilarity is used, the centroid will not correspond to the average of the replicates in the Bray-Curtis space. This is because Bray-Curtis does not satisfy the triangle inequality (Section 5.1.1). A rela-

tionship that can be used to circumvent this problem is the fact that the sum of squared distances from individual points to their group centroid is equal to the sum of squared interpoint distances divided by the number of points in that group (Figure 5.2 B). This implies that if the Bray-Curtis dissimilarities between the data points are known, i.e. the dissimilarity matrix has been calculated, the sum of squared distances to the centroids can easily be calculated without knowing the exact location of the centroid. This way both SS_T and SS_W can be calculated directly from the dissimilarity matrix (Figure 5.3 B and C). From this SS_A can be derived ($SS_A = SS_T - SS_W$) and subsequently the F-statistic can be calculated (Equation 5.2). Often, r^2 values are reported, which indicate the proportion of variance explained by a certain factor (such as time or location in ecological research).

$$F = \frac{SS_A/(a - 1)}{SS_W/(na - a)}. \quad (5.2)$$

To know whether the obtained value for the F-statistic is significantly higher than what would be expected when there would be no difference in the a groups, the distribution for the F-statistic under the null hypothesis has to be known. This distribution can be created by permutating group memberships among the data points. In case the null hypothesis would be true the labels for the rows in the similarity matrix can be shuffled. By calculating F-values, indicated as F' , for all possible random shuffles of the labels the distribution for the F-statistic under the null hypothesis is created. Comparing the obtained F-value with the distribution of the F' -values generates a measure of significance. Considering all possible permutations to calculate F would be computationally intensive, therefore mostly a large, fixed number of permutations is executed (999 in this study).

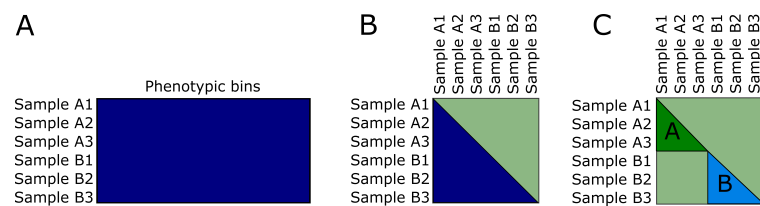


Figure 5.3: Illustration of permanova analysis applied on the phenotypic fingerprints. **A:** The raw data of the phenotypic fingerprints for two samples A en B, each holding three replicates. **B:** The phenotypic fingerprints are converted into a dissimilarity matrix using the Bray-Curtis dissimilarity, which results in a symmetrical dissimilarity matrix. From this dissimilarity matrix SS_T can be calculated by summing the square of all similarities (blue triangle) and dividing by na , the total number of observations. **C:** SS_W can be calculated by summing the square of all similarities within the same group and dividing by n , the number of members in that group. Figure based on Anderson, 2001^[127].

5.1.3 PROcrustean Randomization TEST (PROTEST)

This section is based on Jackson, 1995^[128]. Protest is a randomized test which is based on Procrustes analysis applied on matrices. It compares two matrices by translating, reflecting, rotating and dilating one of the matrices in order to minimize the sum of residuals between the two matrices. In order to assess the significance of the final matrix concordance, a randomized test is used to determine if the sum of residual deviations is less than could be expected to occur by chance.

Consider two matrices for which every point in the first matrix has a corresponding point in the second matrix (Figure 5.4). The fit of these two matrices is maximised by translating, reflecting, rotating and dilating matrix X to some matrix X' . This is done by minimization of the residual sum of squares between the corresponding points in the matrices (i.e. the m^2 statistic, given in Equation 5.3). The final value of m^2 describes how well both matrices fit after the transformation.

$$m_{XY}^2 = \Delta^2(X'_i, Y_i). \quad (5.3)$$

To evaluate the significance of the obtained m^2 value, a permutation test is performed to get an idea about the probability of observing this value. During the permutation, each of the observations from one of the matrices are randomly reshuffled, while maintaining the covariance structure that was present in the matrix.

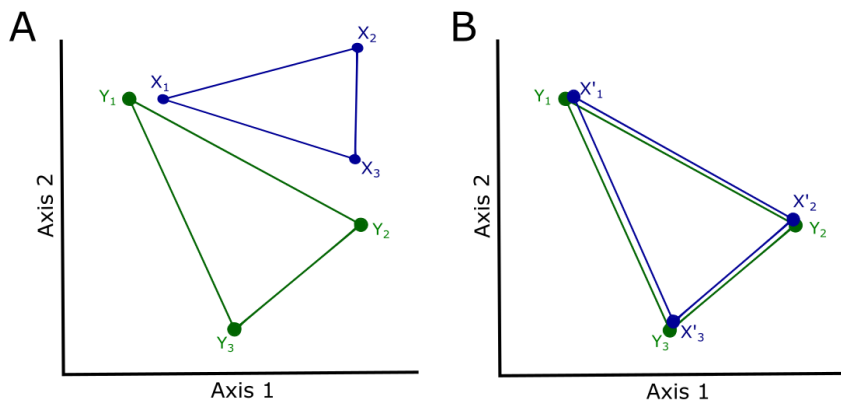


Figure 5.4: Illustration of procrustes analysis. **A:** Consider two matrices X and Y for which every point in the first matrix has a corresponding point in the second matrix. **B:** The fit of these two matrices is maximised by translating, reflecting, rotating and dilating matrix X to some matrix X' . Figure redrafted after Jackson, 1995^[128].

5.1.4 Random forest classifier

A random forest classifier^[129] consists of a set of fully grown decision trees, which are sets of decision boundaries that split up the feature space. The decision tree can be thought of as a set of if- and then-rules via which a new instance can be assigned to a class (Figure 5.5). The new instance will get the label of the class that occurs the most among the data points of the subspace where it is assigned to. Single decision trees tend to ‘overfit’ the data, which means they are good at predicting labels of data which they have seen during the training of the model, but not at predicting labels of unseen data. Random forests avoid this problem by making multiple decision trees. To build this random forest, the dataset is split up into bootstrap samples. For each of the bootstrap samples a decision tree is created. This results in a lot of slightly different solutions to the same problem. Next to the bootstrapping, the algorithm is allowed to choose from only a subset of the predictors at each split. For example, if a random forest is trained on flow cytometry data, at a certain split the algorithm might be allowed to choose from the parameters FL1 and FSC to make its decision boundary, while SSC and FL3 would not be allowed. This approach is used to decorrelate the trees; because in case there is a very strong predictor present in the dataset, most of the trees would use this predictor at the root of the tree, which would result in very similar trees for each of the bootstrap samples. The prediction of all trees is combined into a single prediction by taking the majority vote for all the trees in the forest^[130].

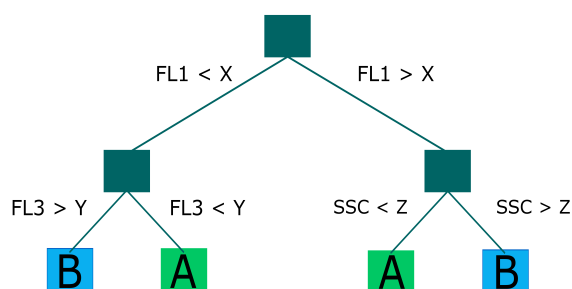


Figure 5.5: Illustration of the concept of a decision tree, applied classification of single-cell FCM data. Through the if- and then-rules an unknown cell can be assigned to species A or species B.

5.1.5 Randomized logistic regression (RLR)

Logistic regression

This section is based on James *et al.*, 2013^[130]. Logistic regression is a linear regression model that is used to model the probability an instance belongs to a certain class ($Y=1$), given the features X of this instance (i.e. the conditional probability $P(Y = 1|X)$, abbreviated as $p(X)$). Since the logistic regression models probabilities, the model output $p(X)$ should range between 0 and 1. Several functions which are suitable for this exist. In logistic regression the logistic function is used (Equation 5.4):

$$p(X) = \frac{\exp(\beta_0 + \sum_{i=1}^k \beta_i x_i)}{1 + \exp(\beta_0 + \sum_{i=1}^k \beta_i x_i)}. \quad (5.4)$$

Linear regression models result in continuous outcomes. Since the outcome of the logistic regression model should be ranging between 0 and 1, the probability $p(X)$ can not be used directly as the model outcome. A solution to this is to rearrange Equation 5.4 to have the ratio of the probabilities on the left side of the equation (Equation 5.5). This ratio of probabilities is called the ‘odds’ and ranges between zero and infinity. The odds is the ratio of the probability the event occurs (class equal to 1), divided by the probability the event does not occur (class equal to 0).

$$\frac{p(X)}{1 - p(X)} = \exp(\beta_0 + \sum_{i=1}^k \beta_i x_i). \quad (5.5)$$

By taking the natural logarithm of Equation 5.5, Equation 5.6 can be obtained. Now the right side of the equation takes the shape of a linear model, the left side is called the ‘logodds’ or ‘logit’.

$$\ln \left(\frac{p(X)}{1 - p(X)} \right) = \beta_0 + \sum_{i=1}^k \beta_i x_i. \quad (5.6)$$

The regression coefficients β of the linear model can be estimated via maximum likelihood estimation.

Regularization

Regularization is a way to reduce complexity of a model and often leads to improved model-predictions. In regularization, a penalty term is added to the objective function that is minimized to estimate the parameters of the model. Different penalties exist, however, only L1-regularization will be discussed here. In a L1-regularization the penalty that is added to the objective function aims to set some of the coefficients to zero. This is known as feature selection and increases the interpretability of the model. The penalty takes the form $\lambda \sum_{i=1}^k |\beta_i|$, with λ the penalization parameter. The intercept coefficient β_0 is generally not penalized. A large value of λ introduces a strong regularization (i.e. a lot of coefficients will be set to zero). The optimal amount of regularization (i.e. optimal value of λ) is generally estimated via cross-validation. The shape of an L1-penalized objective function is given in Equation 5.7.

$$J = \text{objective function} + \lambda \sum_{i=1}^k |\beta_i|. \quad (5.7)$$

Stability selection

Through L1-penalisation a subset of features can be selected. To know whether this subset selection is stable, stability selection can be used^[131]. Stability selection is a combination of subsampling and feature selection. The dataset is split up into n bootstrap samples. For each of these bootstrap samples a feature selection procedure is carried out, which results in n sets of selected features. During each of these feature selections a random subset of coefficients β is perturbed (i.e. the coefficient is scaled with a factor s_i), influencing their chance of being selected (Equation 5.8). By evaluating how many times a variables is selected, a feature ranking can be generated.

$$J = \text{objective function} + \lambda \sum_{i=1}^k \frac{|\beta_i|}{s_i}. \quad (5.8)$$

5.1.6 Receiver operating characteristic (ROC)

The ROC curve gives the true positive and false positive rate for different thresholds a classifier can use to discriminate between two classes. For example, we construct the ROC curves for the classifier which is used for the flow cytometric in silico communities (Figure 5.6). In the ROC curve of species A, a true positive would be a cell which is predicted to be species A and is in reality a member of species A. A false positive would be a cell which is predicted to be species A, but is in reality a member of species B. If a threshold of 0.9 would be applied, this would mean the random forest predicts species A when 90% of the trees in the forest predict species A. At this threshold, the false positive rate would be low (not a lot of cells that are predicted to be species A would in reality be species B), but the true positive rate would also be low since a lot of cells that are in reality species A would not be predicted to be species A with such a stringent threshold. If a threshold of 0.1 would be applied, the true positive rate would be high (a lot of cells that are predicted to be species A, are in reality species A), but also the false positive rate would be high (a lot of cells that are predicted to be species A, are in reality species B). The optimum of the ROC curve is the point (0,1), where there all cells which are predicted to be species A are in reality a member of species A. These ROC curves can be constructed based on the trainingsdata. A random forest classifier normally uses the majority vote of the decision trees to classify a new instance (i.e. threshold of 0.5).

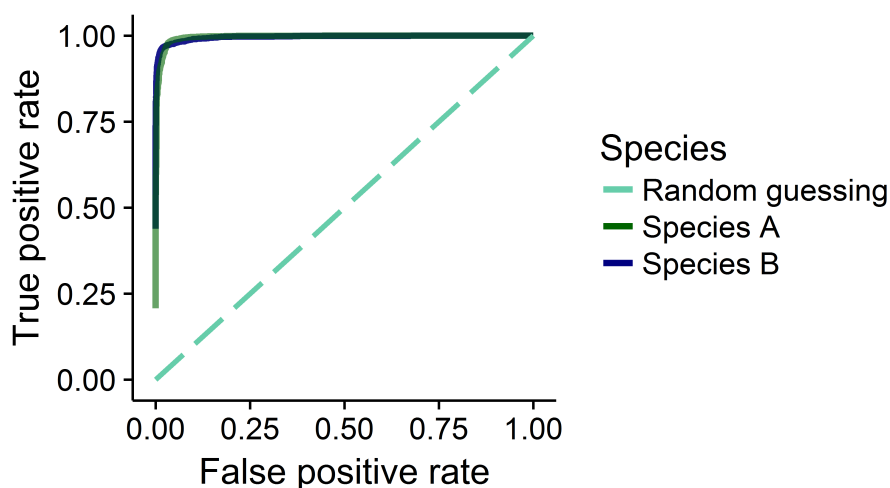


Figure 5.6: Illustration of a ROC curve applied to the flow cytometric in silico communities. The optimum of the curve is located in the upper left corner, at (0,1).

5.1.7 Principal component analysis (PCA)

This section is based on Jolliffe, 2013^[132]. PCA is a technique that is used to visualize multivariate data. It tries to compress high dimensional data into a graph that represents the essence of the information that is captured in the data. The goal is to transfer the data to a 2- or 3D space with the least loss of information, to make it more easily interpretable. This is done by creating new axes in the high-dimensional space that capture as much variance as possible (Figure 5.7). The first axis will be created in the direction that explains the most variance in the high-dimensional space. A second axis will be created orthogonal to the first axis, and in the direction that captures as much as possible of the variance that is still left in the data. New axes will be created until all variance in the data is captured, this implies the number of principal components is lower than or equal to the original number of variables. Correlation is the driver for the principal component analysis since it means some information in the features is redundant. PCA thus turns a large number of correlated variables into an equal or lower number of uncorrelated variables. These new variables are referred to as the ‘principal components’ and are linear combinations of the original variables. To evaluate whether the PCA is doing a good job at presenting the data in the two first components while retaining as much information as possible, the ‘scree plot’ can be evaluated. This plot gives the percentages of variance that are explained by each of the principal components. A good PCA will have large percentages for only a few of the first components and small values for all subsequent components.

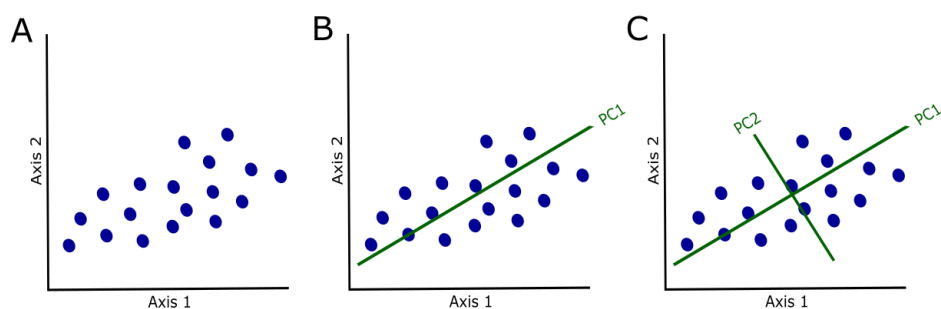


Figure 5.7: Illustration of how PCA finds new axes that explain as much variance as possible. **A:** Consider some data points for which two correlated variables are known. **B:** The first principal component is constructed in the direction that explains the most variance in the dataset. **C:** The second principal component is constructed orthogonal to the first. In this case there is only one option left, however, if the original data would have had more dimensions, multiple options for constructing PC2 would have been possible. In that case, PC2 would have been chosen to capture as much as possible of the variance that was still left in the data.

5.1.8 Principal coordinate analysis (PCoA)

This section is based on Zuur *et al.*, 2007^[133]. Principal coordinate analysis (PCoA) is a kind of metric multidimensional scaling (MDS). It is a technique that is used to visualize dissimilarities. Any dissimilarity measure that is appropriate for the data can be used. The aim of PCoA is to construct a plot in a low-dimensional (2- or 3D) space while respecting the dissimilarities between the data points as well as possible. This way the dissimilarity matrix can easily be interpreted: points that are further away are more distinct, points that are closer to each other are more similar. The exact location of the points is not of interest, it is the spacing from which insight in the underlying behaviour of the data can be gained.

The way this graph is constructed is the following. All dissimilarities between the data points are calculated. The data points are then projected into a high-dimensional Euclidian space respecting all their dissimilarities. By performing a PCA on this high-dimensional space, the points can be represented in a low-dimensional space (for an explanation on PCA see Section 5.1.7). Thus when Euclidean distance is used a PCoA is equivalent to PCA (Section 5.1.7). Note that it is often impossible to map all similarities correctly, the mapping is an approximation of the true similarity matrix^[134].

5.1.9 t-Distributed Stochastic Neighbor Embedding (t-SNE)

This section is based on Van Der Maaten *et al.*, 2008^[135]. t-SNE is a visualisation technique that visualizes the underlying structure of high-dimensional data in a low-dimensional (2- or 3D) space. It specifically aims to preserve the ‘local’ structure of the data.

The first step is constructing pairwise similarities for the high-dimensional data. The distances between the data points in the high-dimensional space are converted into conditional probabilities. The similarity of a data point x_j to another data point x_i is the conditional probability, $p_{j|i}$, that x_i would pick x_j as its neighbor if neighbors were picked in proportion to their probability density under a Gaussian centered at x_i (Figure 5.8). For each point x_i a value of σ_i will be set, so that the conditional probability has a fixed perplexity for each point. This means the bandwidth of the Gaussian in each point is set so that the same number of neighborhood points fall into the Gaussian. This way the Gaussian is adapted to the local densities in the high-dimensional space for each point. The conditional probabilities are then renormalized over all points that contain x_i (Equation 5.9). The final value of similarity p_{ij} between the points x_i and x_j can be calculated as

the average of the conditional probabilities $p_{i|j}$ and $p_{j|i}$.

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma_i^2)}{\sum_{j' \neq i} \exp(-\|x_i - x_{j'}\|^2/2\sigma_i^2)}. \quad (5.9)$$

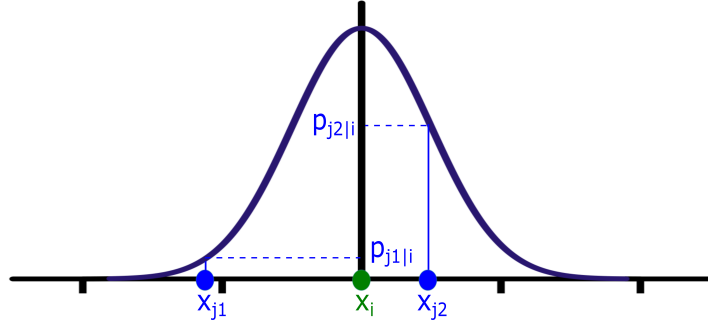


Figure 5.8: The distances between the data points in the high-dimensional space are converted into conditional probabilities. Point x_{j2} is located closer to x_i in the high-dimensional space, compared to x_{j1} . When a Gaussian is centered around x_i , the neighbour point x_{j2} that is located closer will receive a higher conditional probability $p_{j2|i}$ compared to the neighbour point x_{j1} . This way the conditional probability reflects the similarity between the data points in the high-dimensional space.

Now a low-dimensional map (2- or 3D) can be created where each data point in the high-dimensional space will be presented by a point. In a similar way as for the high-dimensional space, conditional probabilities q_{ij} for the points in this low-dimensional space can be calculated. However, in this low-dimensional space a Student-t distribution is used instead of a Gaussian distribution to calculate the conditional probabilities (see further). The aim is to have the similarities in the low-dimensional map representing the similarities in the high-dimensional space. In other words we want q_{ij} to be representing p_{ij} as well as possible. If q_{ij} 's and p_{ij} 's are similar, the original structure of the data is preserved well. Therefore an objective function that will measure the discrepancy between similarities in the high-dimensional space and similarities in the low-dimensional map is needed. For this a cost function C based on the Kullback Leibler divergence is applied (Equation 5.10). By minimizing this cost function, the discrepancy between p_{ij} and q_{ij} is minimized while giving most importance to preserving the local structure of the data. This is due to the fact that Kullback Leibler divergence is asymmetric. When two points which are far away from each other in the high-dimensional space (small p_{ij}) are close to each other in the low-dimensional space (high q_{ij}), this will lead to a smaller penalty compared to the situation where two points which are close to each other in the high-dimensional space (high p_{ij}) are far away from each other in the

low-dimensional space (small q_{ij}). This way more importance is given to the local data structure.

$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \left(\frac{p_{ij}}{q_{ij}} \right). \quad (5.10)$$

The reason for the use of a Student-t distribution when calculating the conditional probabilities in the low-dimensional space is that this distribution has a more heavy ‘tail’ compared to a Gaussian distribution. When high-dimensional data is mapped down to a lower number of dimensions, it is impossible to preserve all similarities of the high-dimensional space exactly. t-SNE gives more importance to the local structure of the data which results in the fact that data points that are far away in the high-dimensional space often get mapped too far away in the low-dimensional space. This will cause a small contribution to the cost function. However, in these large datasets there are often many points that are mapped to far away from each other. All these (small) contributions to the cost function can cause the low-dimensional graph to collapse and thus no longer represent the clusters that are present in the data. Using a heavy-tailed distribution explicitly allows data points that are far away in the high-dimensional space to get mapped far away in the low-dimensional space, avoiding the accumulation of all these small contributions to the cost function.

t-SNE is not designed to conserve the global structure in the data, however, when one wants to mitigate this issue, the t-SNE algorithm can be seeded with the result of a PCA. This way the local structure will be embedded in a graph that is already organized in a way to respect the more global structure in the data. The reason why t-SNE sometimes does a better job than PCA at visualizing the structure of the data is that PCA focuses on explaining as much variance as possible in the dataset and is therefore focusing at placing observations that are far away in the high-dimensional space also far away in the low-dimensional map. Observations that are present in between are simply projected onto the principal component without taking their interrelationships into account (Figure 5.9).

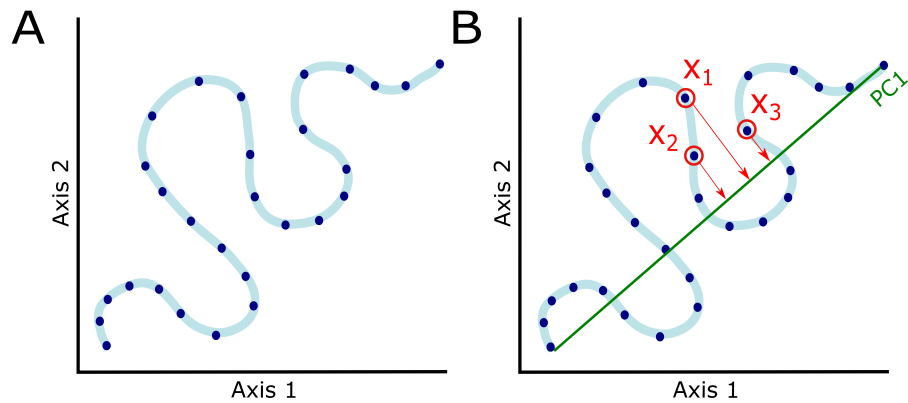


Figure 5.9: Illustration of a situation where t-SNE might lead to a better visualisation of the underlying structure of high-dimensional data compared to PCA. **A:** Consider data that is organized according to some structure (indicated with the light blue line) in the high-dimensional space. **B:** If a PCA would be performed on this dataset not all relationships between data points would be respected correctly. The axis that explains the most variance would be chosen as PC1 (indicated with the green line). The location of the points x_1 , x_2 and x_3 on the principal component would not be reflecting their true relationship. The point x_1 would be situated in between x_2 and x_3 . While, if you would take into account the local structure of the data x_2 should be in between x_1 and x_3 , with x_1 more similar to x_2 compared to x_3 .

5.2 Supplementary figures

5.2.1 Raman spectra nucleic acids

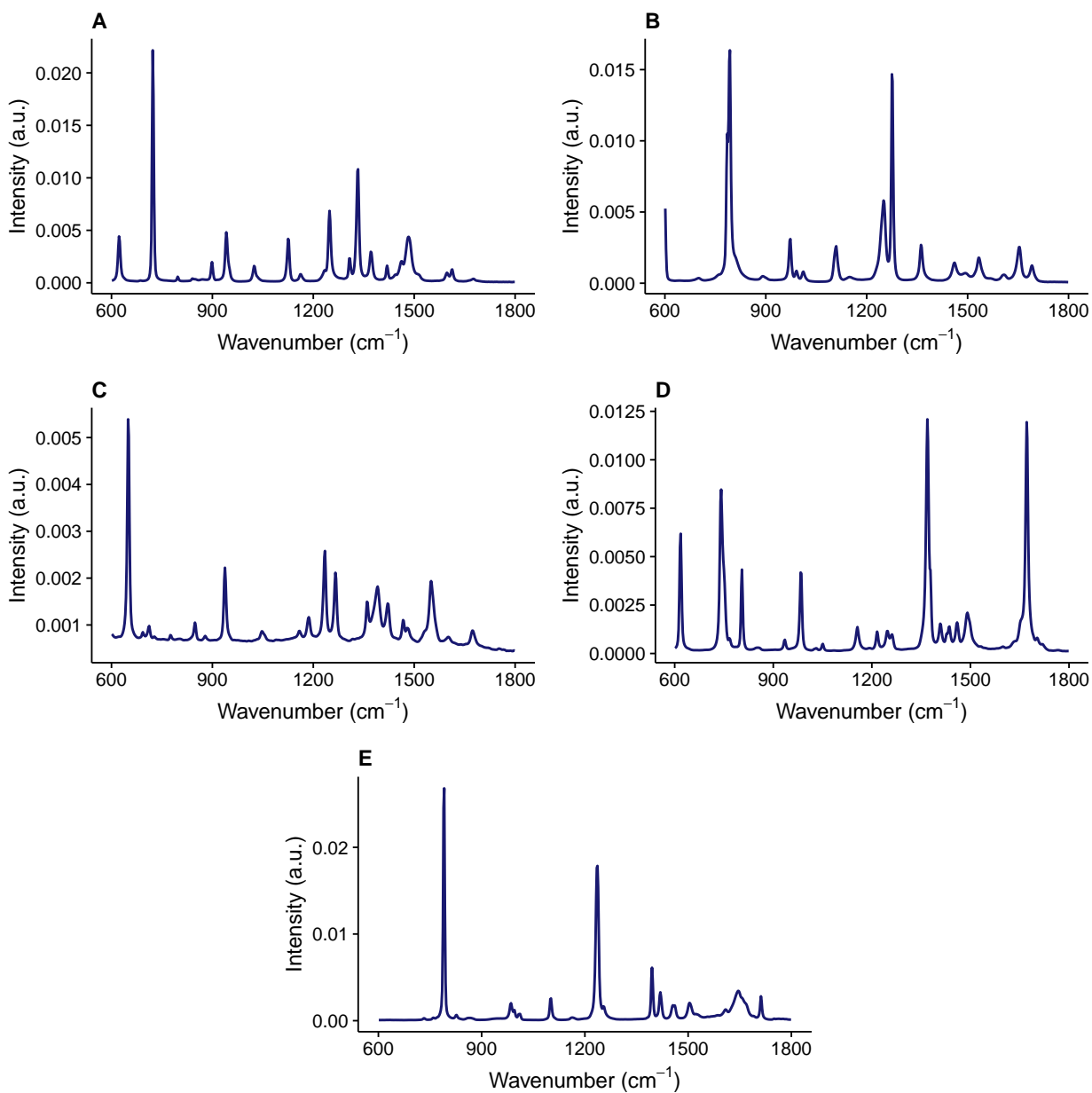


Figure 5.10: Raman spectra of DNA and RNA bases: adenine (A), cytosine (B), guanine (C), thymine (D) and uracil (E). Preprocessed data was retrieved from the study of De Gelder *et al.*, 2007^[77].

5.2.2 Lake isolates

Sequences of the *Limnohabitans* lake isolates were obtained from Lake Michigan during the 2016 Summer survey. In order to situate the available isolates within the genus of *Limnohabitans*, a phylogenetic tree was created. The isolate sequences were aligned to the Silva database (v123) using Mothur (v1.38, seed = 777). Fasttree (v2.1.10) was used to create the phylogenetic tree. The tree was visualised using iTOL (v3.4.1). The result is shown in Figure 5.11. All three lake isolates belong to the Lhab-A1 clade.

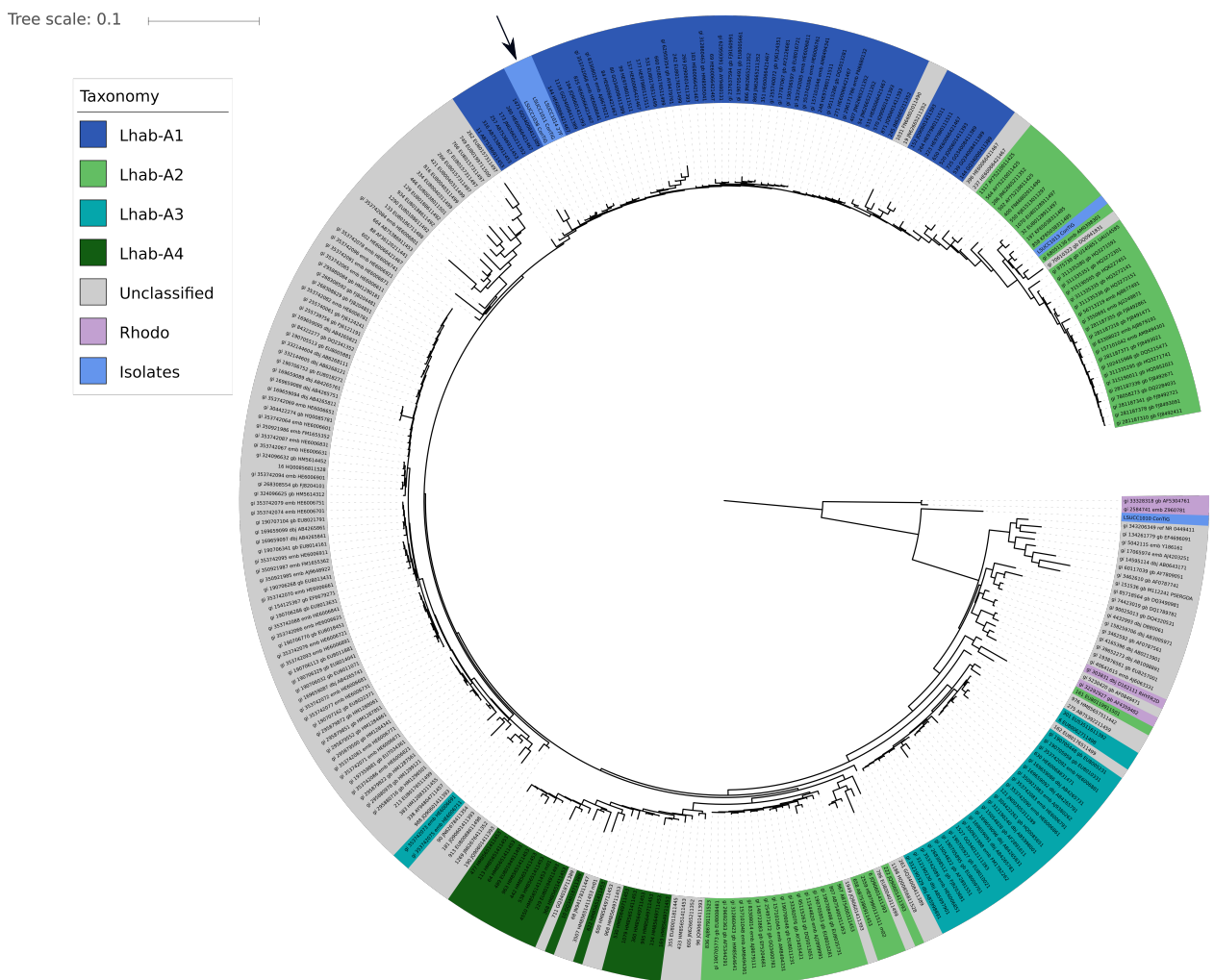


Figure 5.11: Phylogenetic tree to situate the available lake isolates within the genus of *Limnohabitans*. The three isolates are indicated with an arrow. All three belong to the Lhab-A1 clade.

5.3 Supplementary tables

5.3.1 Predicted relative abundances experiment 1

Table 5.1: Predicted relative abundances for both community members in the mixed cultures during experiment 1. Predictions were made for each of the biological replicates ($n = 3$) separately. The first column, ‘Approach’, corresponds to the identifiers that were used in Figure 3.6 (A = fingerprints of the axenic cultures at 24h were used as training data, B = fingerprints of the axenic cultures at the corresponding time point were used as training data, C = fingerprints of the cocultures at the corresponding time point were used as training data).

Approach	Time (h)	Replicate	Abundance A (%)	Abundance B (%)	Accuracy test set (%)
A	24	1	70.19	29.81	99.70
		2	72.50	27.50	
		3	65.00	35.00	
	48	1	0.42	99.58	
		2	0.25	99.75	
		3	0.29	99.71	
	72	1	0.24	99.76	
		2	0.25	99.75	
		3	0.33	99.67	
B	24	1	70.19	29.81	99.70
		2	72.50	27.50	
		3	65.00	35.00	
	48	1	0.59	99.41	99.95
		2	0.61	99.39	
		3	0.73	99.27	
	72	1	0.11	99.89	99.97
		2	0.22	99.78	
		3	0.22	99.78	
C	24	1	62.23	37.77	99.7
		2	63.93	37.07	
		3	63.40	36.60	
	48	1	32.14	67.86	98.83
		2	31.83	68.17	
		3	34.98	65.02	
	72	1	29.47	70.83	97.20
		2	28.58	71.42	
		3	29.45	70.55	

5.3.2 Criteria for cell classification

Table 5.2: Optimal thresholds that were used as decision boundary when creating the in silico communities based on the axenic cultures. These thresholds were determined as the point closest to (0,1) on the ROC curve when the random forest was trained on data of the axenic cultures. (TP = true positive, FP = false positive)

Time (h)	Threshold A	TP A	FP A	Threshold B	TP B	FP B
24	0.72	0.9983	0.0037	0.30	0.9963	0.0017
48	0.65	0.9993	0.0000	0.44	1.0000	0.0017
72	0.64	0.9990	0.0000	0.57	1.0000	0.0017

Table 5.3: Optimal thresholds that were used as decision boundary when creating the in silico communities based on the cocultures. These thresholds were determined as the point closest to (0,1) on the ROC curve when the random forest was trained on data of the cocultures. (TP = true positive, FP = false positive)

Time (h)	Threshold A	TP A	FP A	Threshold B	TP B	FP B
24	0.59	0.9841	0.0294	0.42	0.9705	0.0159
48	0.63	0.9924	0.0147	0.38	0.9853	0.0076
72	0.59	0.9751	0.0304	0.42	0.9696	0.0249

5.3.3 Relative abundances experiment 3

Table 5.4: Relative abundances of the two species A and B in as determined in the mixed community with a *gfp*-labeled strain of species A.

Time (h)	Replicate	Relative abundance A (%)	Relative abundance B (%)
24	1	76.55	23.45
	2	82.08	17.92
	3	80.36	19.64
48	1	72.54	27.46
	2	60.43	39.57
	3	69.01	30.99
72	1	54.30	45.67
	2	56.27	43.73
	3	66.63	33.37