

2012

# Analysis of Behavior in Populations of Swimming Microbes

David Jordan

Follow this and additional works at: [http://digitalcommons.rockefeller.edu/student\\_theses\\_and\\_dissertations](http://digitalcommons.rockefeller.edu/student_theses_and_dissertations)

 Part of the [Life Sciences Commons](#)

---

## Recommended Citation

Jordan, David, "Analysis of Behavior in Populations of Swimming Microbes" (2012). *Student Theses and Dissertations*. Paper 166.

This Thesis is brought to you for free and open access by Digital Commons @ RU. It has been accepted for inclusion in Student Theses and Dissertations by an authorized administrator of Digital Commons @ RU. For more information, please contact [mcsweej@mail.rockefeller.edu](mailto:mcsweej@mail.rockefeller.edu).



# ANALYSIS OF BEHAVIOR IN POPULATIONS OF SWIMMING MICROBES

A Thesis Presented to the Faculty of

The Rockefeller University

in Partial Fulfillment of the Requirements for

the degree of Doctor of Philosophy

by

David Jordan

June 2012



# ANALYSIS OF BEHAVIOR IN POPULATIONS OF SWIMMING MICROBES

David Jordan, Ph.D.  
The Rockefeller University 2012

This work describes our work developing an experimental biological system to study patterns of behavioral variability. We selected motility as the behavior of interest because it is common throughout biology and can be recorded and analyzed relatively easily. We chose to work with a microbe, *Tetrahymena thermophila*, as a model organism for these studies; it is easy to grow in the laboratory in controlled conditions, has a relatively short generation time, and is large enough to for its motions to be easily imaged. To achieve the imaging, we developed a set of low cost digital video microscopes. Concurrently, we wrote custom software to create trajectories of movements from the recorded movies. Consumer webcams provided high temporal and spatial resolution at low cost, and custom microfluidic devices allowed organisms to be isolated and studied in a well-controlled environment. Simultaneous tracking of multiple individuals, while retaining the identity of each, allowed experiments to span multiple generations. Further, we developed a method of characterizing the swimming behaviors using histograms of linear and angular speeds, which did not rely on explicit modeling or scoring of stereotyped behaviors, and used it to quantitatively measure the similarity between behaviors. These similarities were computed using a relative entropy based metric called the Jensen Shannon divergence. Using this framework, we measured patterns of behavioral changes, both within individual lifetimes and between different individuals in a population. These changes were quantified over time

scales that ranged from minutes to hours and even between generations. We measured all of these in a variety of environments, and catalogued the effects of changing the environment. We used the similarity measurements generated from the above analysis to generate a low-(two-) dimensional representation of the behaviors, which led to convenient visualization of the patterns of behavioral change and variability. In addition, we performed experiments using artificial selection that provided evidence that this low dimensional representation may be of biological relevance.

# Acknowledgments

The time I have spent in the lab, on the work presented here, has been wonderful, unique, stimulating, and wholly pleasurable. I would like to thank the people that made it so,

First and foremost, I would like to thank my advisor, Dr. Stanislas Leibler. Working with Stan has been a privilege and an honor. It is hard to describe the unique mix of enthusiasm, excitement, integrity, and fatalism that makes Stan such an inspiration to work with.

I also had the good fortune to work with Seppe Kuehn, whose patience and pragmatism, combined with his work ethic and easygoing personality have made the experience a singular pleasure. I am grateful to him for his willingness to always share his knowledge and for often reminding me to trust in the data.

I would like also to thank Eleni Katifori. In the early days of this project, we had many fun and enlightening forays into data analysis and she was instrumental in the development of the method we finally settled on.

My committee, Dr. James Hudspeth and Dr. Albert Libchaber. For inspiring me to be rigorous and to always ask myself "so what does this all mean?" I am grateful for the guidance and insight provided both at our annual meetings, and during discussions we had in the interim.

David Huse, for serving as the external examiner for my defense, and for stimulating discussions that continue to influence how I think about these data.

Finally I would like to thank the Dean's Office and the Rockefeller Graduate Program for truly allowing me to focus on the science.

# Table of Contents

Table of Contents	iv
List of Figures	vi
List of Tables	viii
<b>1 Introduction to Behavioral Analysis in Model Systems</b>	<b>1</b>
1.1 Behavior in Biological Systems . . . . .	2
1.2 Motility as a Behavioral Model . . . . .	5
1.3 System Design . . . . .	7
1.4 Motivations . . . . .	8
<b>2 Digital Microscopy for Tracking Swimming Microbes</b>	<b>13</b>
2.1 Hardware . . . . .	14
2.1.1 Imaging System . . . . .	14
2.1.2 PDMS Chambers . . . . .	15
2.1.3 Temperature Control . . . . .	19
2.2 Wetware . . . . .	19
2.2.1 <i>Tetrahymena sp</i> Biology . . . . .	19
2.2.2 Growth Media . . . . .	22
2.3 Software . . . . .	23
2.3.1 Image Analysis . . . . .	23
2.3.2 Tracking . . . . .	24
<b>3 Analysis of Behavioral Patterns in <i>T. thermophila</i></b>	<b>27</b>
3.1 Preliminary Observations . . . . .	27
3.2 Transformation of Data . . . . .	31
3.3 Quantitative Representation of Behaviors . . . . .	31
3.3.1 Measuring Statistical Distance . . . . .	33
3.4 Patterns of Individual Behavioral Change . . . . .	35
3.4.1 Changeability . . . . .	36
3.4.2 Behavioral Memory . . . . .	38
3.5 Behavioral Variability in Populations . . . . .	40
3.5.1 Individuality . . . . .	41
3.5.2 Plasticity . . . . .	43

3.5.3	Heritability . . . . .	43
3.6	Dimensionality Reduction . . . . .	47
3.6.1	Multidimensional Scaling . . . . .	48
3.6.2	Evaluating Embeddings . . . . .	49
3.6.3	Trajectories in Behavioral Space . . . . .	51
3.6.4	Effects of Environmental and Genetic Perturbations . . . . .	52
3.7	Dimensions in Behavioral Space . . . . .	55
3.7.1	Perceptual Mapping . . . . .	55
3.7.2	Evidence from Selection in Behavioral Space . . . . .	56
<b>4</b>	<b>Conclusions and Future Directions</b>	<b>64</b>
4.1	Conclusions . . . . .	64
4.2	Future Directions . . . . .	65
<b>A</b>	<b>Chamber Homogeneity and Boundary Effects</b>	<b>68</b>
A.1	Chamber Depth . . . . .	69
A.2	Walls . . . . .	69
A.3	Chamber Isotropy . . . . .	76
A.4	Nutrient Depletion . . . . .	78
<b>B</b>	<b>System Details</b>	<b>79</b>
B.1	Temperature Control . . . . .	79
B.2	Chamber Fabrication . . . . .	79
B.3	Detection . . . . .	81
B.4	Cost Functions . . . . .	82
<b>C</b>	<b>Behavioral Analysis</b>	<b>86</b>
C.1	Comparing behaviors through distributions divergence . . . . .	87
C.2	Changeability: variations of behavior during lifetime . . . . .	89
C.3	Individuality: Differences of behavior among individuals . . . . .	92
<b>D</b>	<b>Tracking, Error and Uncertainty</b>	<b>94</b>
D.1	Tracking Fidelity . . . . .	94
D.2	Imaging Uncertainty . . . . .	95
D.3	Divergence Estimators . . . . .	95
D.3.1	Bias . . . . .	96
D.3.2	Uncertainty . . . . .	97
<b>E</b>	<b>Changeability and Individuality Distributions</b>	<b>98</b>
<b>F</b>	<b>Protocols</b>	<b>101</b>
F.1	<i>Tetrahymena</i> sp. Information . . . . .	101
F.2	Cell Culture . . . . .	101
	<b>Bibliography</b>	<b>103</b>



# List of Figures

2.1	Schematic and Photograph of Imaging System . . . . .	16
2.2	Schematic of PDMS Microfluidic Chamber . . . . .	18
2.3	Scanning Electron Micrograph of <i>T. thermophila</i> . . . . .	22
2.4	Example Image from Digital Microscope . . . . .	25
3.1	Distributions of Generation Times in 5 Environments . . . . .	29
3.2	Generation Times of Sisters are Perfectly Correlated with Post Division Size . . . . .	30
3.3	Example Trajectory for a Single Individual . . . . .	32
3.4	Changeability Matrices Show Patterns of Individual Behavioral Change	37
3.5	Behavioral Variability in Populations . . . . .	46
3.6	Scree Plot of Embeddings for Different Populations . . . . .	50
3.7	Behavioral Trajectories from Multidimensional Scaling of Changeabil- ity Matrices . . . . .	52
3.8	Raw Embedding and Contours of Kernel Density Estimate . . . . .	53
3.9	Environmental and Genetic Perturbations are Apparent in Behavioral Space Embeddings . . . . .	54
3.10	Perceptual map of a single trajectory in 2 dimensions . . . . .	56
3.11	Changes along the red axis correlate with changes in high-speed mode density . . . . .	57
3.12	Changes along of the red axis show little change in high-speed mode location . . . . .	58
3.13	Changes along the blue axis correlate with changes in high-speed mode location . . . . .	59
3.14	Longer lived sisters after divisions always have a high-speed mode that is faster . . . . .	61
3.15	Selection Induced Behavioral Changes can be Observed in Behavioral Space . . . . .	62
A.1	Averaged Histograms of Swimming Speed for Different Chamber Depths	70
A.2	Ratio of Swimming Speeds in Proximity to a Wall as a fraction of Speed Far From the Wall . . . . .	71
A.3	Distributions of JS divergence between Populations Assayed in Cham- bers of Different Depths . . . . .	72
A.4	Radial Distribution Function Near the Wall . . . . .	73

A.5	Distribution of Durations Spent Near the Wall . . . . .	74
A.6	Fraction of Lifetime Spent Near the Wall . . . . .	75
A.7	Behavioral Similarity for Behaviors Punctuated by an Interaction with the Wall . . . . .	76
A.8	Spatial Correlation Between Individual Distribution Functions . . . .	77
B.1	Circuit Diagram for Thermistor Amplifier . . . . .	80
B.2	Circuit Diagram for Peltier Driver . . . . .	80
C.1	Histograms of Actions for a Simple Example . . . . .	88
C.2	Measuring Individual Behavioral Changes . . . . .	91
C.3	Measuring Differences Between Individuals . . . . .	93
E.1	Changeability Distributions for All Environments for <i>T. thermophila</i>	99
E.2	Individuality Distributions for All Environments for <i>T. thermophila</i> .	100

# List of Tables

2.1	Abbreviations for Growth Media . . . . .	23
3.1	Quantities which describe behavioral dynamics . . . . .	47
F.1	<i>Tetrahymena sp.</i> Strain Information . . . . .	101

# Chapter 1

## Introduction to Behavioral Analysis in Model Systems

In general, behavior refers to how a system changes in response to perturbations. When we say a function behaves non-linearly, we are really saying that when you vary its input, the output changes in a non-proportional way. In a physical sense, the parameters that vary are often measurable quantities. For example, when we say a substance behaves like a gas, we are referring to its tendency to expand to the shape of its container. We can quantify this tendency by measuring its change in pressure as we perturb its volume (change the size of its container). Furthermore, we can summarize the relationships between these important parameters mathematically, e.g. Boyle's law. We can also include other important parameters to obtain more complete descriptions with fewer assumptions, such as the incorporation of temperature and particle number into the Ideal Gas Law. The Ideal Gas Law is a description, using only a few meaningful variables, which subsumes all of the molecular details of the dynamics of the individual gas molecules.

Simplifying descriptions for complex non-equilibrium systems are not generally known. In particular, in biology, whether there are measurable coarse-grained quan-

tities which will yield meaningful reduced complexity descriptions, and what those quantities might be, are still open questions. While these questions in general remain unanswered, there is evidence that such descriptions might be possible. For example, empirical quantities that relate the growth rate of bacteria to some properties of cell composition have been summarized in a set of bacterial growth laws [41] and effective parameters which relate average motility and density have been used to describe pattern formation in growing bacterial colonies [11]. This work presents our initial investigations in to how we might look for evidence of the existence of such quantities for another behavioral model. This section will describe what we believe to be the necessary requirements if such an experimental system is to be successful.

## 1.1 Behavior in Biological Systems

In biology, the ability of an organism respond appropriately to changes in its environment will determine its evolutionary success. Biological systems exhibit adaptive behaviors. Adaptation is the notion that among the variables that describe the state of the system, some are able to change rapidly in response to environmental change, while others can only change slowly. If this is the case, and we can make a reasonable separation of time scales, the fast dynamics of the system will see the slowly changing variables as parameters. The changing of these slowly varying parameters is termed adaptation [21].

Time scale separation may be difficult, however, as responses may depend on multiple environmental parameters that change on different time scales. To add

further complexity, responses which depend on the state of the environment may also feedback and change the environment. Thus, both the responses themselves, and their effects on the environment, may have many characteristic time scales. For example, in response to the presence of a particular sugar, a bacterium may begin the process of making cellular machinery to utilize that sugar in metabolism within minutes [40], however, these systems may remain induced for hours, even if the sugar concentration changes [49]. The response feeds back to the environment as the sugar is metabolized by the newly made enzymes and its concentration depleted. While these processes are occurring, other environmental parameters, such as temperature, might fluctuate on time scales of seconds or minutes. Furthermore, over long times, a new mutation may arise that allows for better utilization of that sugar, a response that has time a time scale of very many generations. Because the important time scales are unknown, the ideal system would record the behavior at high time resolution, for long durations, so that the appropriate time scales could be determined from the data.

Because the slowly changing variables are unknown, long duration behaviors that were initiated prior to the beginning of our experiment may have measurable effects. This phenomenon is known in *Escherichia coli*, where the response of the organism to temperature gradients is dependent on a persistent state associated with the growth phase of the population from which it came. Thus if you take cells from early phases and late phases of growth, they will have different thermotactic responses even when assayed in the same conditions [39]. We do not know all such slowly changing parameters, nor do we know the duration of measurements required to discover them. If they vary within a generation, we can measure their decorrelation time with a quan-

tity we call behavioral memory, if they are persistent across generations, this will be apparent as behavioral heritability, similar to the persistence of the induction of the *lac* operon described in [49].

In addition, if the future is uncertain, behavioral responses may be probabilistic. Such responses require a statistical description. Probabilistic responses among individuals in a population is a form of phenotypic variation and if that population is isogenic, it is called non-genetic individuality [43], and has been observed in *E. coli*. When variation is employed to cope with uncertainty and fluctuations in the environment, it is sometimes called bet-hedging [3]. Because of this variability, the ideal system would allow for many measurements to be made simultaneously, allowing for the collection of well sampled, population-level statistics. With such statistics we could measure the variety of behaviors presented by different individuals, a measure of the individuality of the population.

Organisms are not likely to have swimming behaviors that are adapted to laboratory conditions, which have generally been chosen to give optimal growth. Over evolutionary time scales, extant organisms likely have adapted to the statistics and to the temporal and spatial correlations of their natural environments. Because we do not know these statistics, or how they have changed over the evolutionary history of the organism we would like system where many environments can be measured simultaneously and where environments can be varied in many dimensions in a controlled way.

Experimental studies of behavior are often carried out in animals and microbes. Laboratory studies have elucidated genetic aspects of circadian clocks in the fruit fly

*Drosophila melanogaster* [26], environmental and genetic contributions to foraging strategies in the nematode *Caenorhabditis elegans* [4], and the molecular details of chemotaxis in the bacteria *E. coli* [1]. These model organisms have been chosen for experimental studies because they are easy to maintain in a laboratory setting, have reasonably short generation times, and are amenable to controlled genetic and environmental perturbations. Well characterized behavioral traits, which generally require only short term measurements, can be studied easily in the laboratory in the model systems described above. However, long term measurements of the entire behavioral repertoire of a population of many individuals and their progeny has not been reported. Model systems offer distinct advantages over studying behavior in the field. The experimenter can readily control many aspects of the environment and often has available a variety genetic tools. In particular, microbial systems, in contrast to animal models such as *D. melanogaster* and *C. elegans*, are well suited to the sorts of studies we were interested in undertaking due to their short generation times (hours). This allows many generations to be recorded in a reasonable amount of time, and for small numbers of individuals to be quickly expanded into very large ( $10^6$ ) populations.

## 1.2 Motility as a Behavioral Model

While all the responses of an organism constitute its behaviors, including metabolism and reproduction, one of the most conspicuous behaviors animals undertake is movement. Motion is a common feature of biological systems [14] and provides an ad-



vantage to organisms that live in environments which are spatially heterogeneous on length scales longer than the size of the organism. Although the outputs of metabolism require sophisticated equipment to measure, motion can be measured relatively simply using digital imaging. If we restrict our interest to motions that an organism can undertake, a catalogue of such actions can be obtained by recording such motions and generating trajectories.

One of the earliest model systems for tracking was that of *E. coli* in the lab of Howard Berg [6]. This system used a three-axis motorized stage with feedback from the imaging system to keep a single bacterium in focus in the field of view of the microscope. We wanted to avoid a system that relied on mechanical motors, which can be costly and unreliable in long term experiments. In addition, this methodology allowed only a single organism to be tracked, while we wanted to track many organisms simultaneously. This scheme has been parallelized allowing the recording of trajectories for up to half a generation time for many *E. coli*, but these were not in the same arena during tracking [5]. The current state of the art in long term experiments relies on trapping organisms for imaging either optically [32], or mechanically [22] and is still limited to one individual at a time. Tracking of multiple individuals in the same arena, up to thirty simultaneously [9], has been achieved, but only for short periods of time, less than 0.1% of the lifespan. Our goal, to track multiple freely moving organisms in the same arena, for many generations, has not, to our knowledge, been achieved.

Finally, the ideal methodology would include a quantitative characterization of recorded motions, a simplified description of the actions that an organism presents

over its lifetime. Traditionally, behavioral characterizations are made by scoring, which breaks sequences of actions into pre-recognized classes called stereotypes [9]. Examples of stereotyped behaviors might include walking and running. Although scoring allows the measurement of some differences, for example, by computing the difference between how often individuals run, it breaks down when novel behaviors emerge or if the boundaries between stereotypes become unclear. In addition it does not allow stereotyped behaviors to be compared to one another. The ideal characterization should avoid scoring and provide an explicit, quantitative measure of behavioral differences that allows any behavior to be compared to any other.

### **1.3 System Design**

In summary, we will describe the development of an experimental system to measure behavior in a biological system. Because motility is common in biology, important for survival, and has an output which is fairly easy to measure, it was chosen as a behavioral model. Long-duration, (ideally many generations) high-time-resolution experiments allowed the measurement of behavioral changes on many time scales simultaneously as well as the observation of long lasting behavioral correlations. Construction of multiple imaging systems made it feasible to measure behavior in populations, allowing us to generate statistical descriptions to accommodate probabilistic responses. Furthermore, replicate systems also allowed us to sample a variety of environments. These experiments were carried out in a microbe-based model system, for microbes can be grown easily in well controlled conditions and have relatively short

generation times. We recorded swimming motions of organisms confined to a quasi two-dimensional chamber using video microscopes, and generated characterizations of the motion that avoided scoring and that allowed for quantitative comparisons between behaviors, even on different time scales.

## 1.4 Motivations

In even the simplest model of evolution, phenotypic variation is of central importance. Phenotypic variation is the raw material of natural selection, and is the product of complex ecological interactions between organisms of different evolutionary histories.

Organisms are not generated *de novo*, but arise from the replication of existing organisms. As such, all extant organisms are the realizations of historical contingencies that extend all the way back to the origin of life. The evolutionary history of two organisms, which includes selection upon past genetic variation, can lead to different phenotypes, even in the same current conditions. The dependence between the current phenotype and selection for past genetic variants is often referred to as the “genotype phenotype map”. There are many other mechanisms of generating heritable phenotypic variation that do not rely on the genetic mechanism. For example, the phenotypic variety of differentiated cells in a multicellular organism depends strongly on each cell’s life history, in particular, on its environment during development. In general, the non-genetic mechanisms are harder to study, leading to an emphasis on genetics.

The preoccupation with the genetic mechanism has driven the development of

technologies capable of measuring genetic variation directly. The development of methodologies to measure phenotypic variation have lagged behind. Many phenotypic traits, particularly those related to structure and development, do not vary much during the lifetime of an individual and are fairly easy to characterize. One such example, the number of bristles on the abdomen of *Drosophila pseudoobscura* [19], is easily characterized by a single number and does not change during the lifetime of an organism. However, other aspects of phenotypes, in particular behavioral ones, are much harder to characterize and can change significantly during an organisms lifetime. If we consider motility as a behavioral phenotype, it is clear that the variability in this phenotype within an organism's lifetime is significant, as motions are known to change in response to environmental stimuli of time-scales much shorter than the typical lifespan of the organism. What kinds of measurements are needed to characterize such a phenotype?

Because motile behaviors are difficult to characterize in general, most studies of motility limit themselves to measuring taxis, that is, motion in a gradient. Many taxis responses can be easily characterized by short measurements before and after the application of a stimulus [44], [43]. However, it is known that even in a homogeneous environment, organisms undergo complex motions and do not simply perform random walks [6], [27]. In fact, Korobkova and colleagues showed that motile behaviors can show temporal fluctuations due to molecular noise in a single bacterium. Because of this intrinsic variation, if we seek to characterize behavior, even for a single individual in homogeneous environment, short measurements are unlikely to be sufficient. Ideally, we would like to measure behaviors for the entire lifetime of the

individual.

Although recently full lifetime behavioral measurements [22], [32], have been reported, these reports present data from individuals that are help in place either mechanically or optically. The majority of studies of behavior, and all of those done in freely moving organisms, have been done with sub-lifetime measurements. The longest measurements of freely moving organisms that have been reported can measure individuals for up to 1/40th of a lifetime, in *E. coli*. These ratios are even lower for higher organisms such as *C. elegans* (1/2000) [16], and for *Drosophila* (1/4000) [51]. As a comparison, the longest measurements currently available would be equivalent to characterizing a lifetime of human behaviors with a measurement of about 2 years. With this in mind, we sought to measure how well short observations represent full lifetime behaviors. To to this, we needed to make the measurements of full lifetime behaviors in freely moving organisms.

Once we have a measurement of full-lifetime behaviors in a individual, we would like to replicate that measurement in many individuals in a population. With measurements in multiple individuals, we could compare within individual live-time variability to the variability between individuals in a population. Such variability in isogenic populations in *E. coli* was shown originally by Spudich and Koshland [43], who called it non-genetic individuality. However, because they did not have full lifetime measurements, to was not clear that the individuality observed was not due in some part to within individual variability. If behaviors can vary among isogenic individuals in identical environments, it is important to characterize how much variability is expected. This measurement will establish a baseline which can be used

to determine how many individuals must be measured to characterize a population. Furthermore, when studying the effect of an experimental perturbation on behavior, these measurements will allow us to determine whether differences that are observed are significant relative to the differences we would expect due to individuality.

With characterizations of the behavioral variability of individuals in populations, we might wonder if that variability is characteristic of natural populations, or is a result of our particular choice of growth environment. The natural environments of most model organisms are not well characterized, with some exceptions [13], which makes it difficult to address this question directly with experiments. However, we can ask a related question, which is, when populations are grown in different environments, do populations exhibit the same phenotypes. The phenomenon of different phenotypes arising in different environments is well known, and is called phenotypic plasticity. Both behavioral [20] and developmental [19] phenotypes have been shown to be plastic in model organisms. One advantage of having full life-time measurements of behavior from multiple individuals is that the variability in behavior in individuals or in populations can be compared to the variability between populations grown in different environments.

Finally, for a selection process to become a process of adaptive evolution, the mechanisms which generate successful phenotypic variation must be inherited from one generation to the next. Because these mechanisms are not simple genetic, instead of measuring behavioral differences between different genetic backgrounds, we would like to measure correlations between generations directly. This allows for a measurement of heritability that is based directly on generational distance, and does

not rely on a particular model of the inheritance mechanism. To do this requires not only measurements of full-lifetimes, but also measurements of many generations, with accurate phylogenies. To achieve this, organisms must be confined, and not allowed to hide or escape during experiments. In addition, we must be able to image them frequently enough such that identity is maintained.

The goal of these measurements will be to answer the following questions. First within an individual's lifetime, how well do short measurements represent the full lifetime. This question allows us to measure intrinsic variability in a single individual, and can be answered by making full lifetime measurements and we present this measurement as a quantity called changeability. With full lifetime measurements, we can next ask how well does the observation of one individual represent the behaviors of a population? We call this measurement individuality. Estimates of individuality will allow us to establish a base line against which to determine whether experimental treatments change behavior of populations. One of the first experimental treatments we will try to change behavior in populations will be altering the environment. Determining whether environmental differences generate behavioral differences will require the measurement of populations in a variety of environments. With such measurements, we can ask whether the behaviors presented in one environment are representative of those observed in any other environment. Lastly, we will present a measurement about whether the similarity in behaviors is correlated with relatedness.

## Chapter 2

# Digital Microscopy for Tracking Swimming Microbes

The system described at the end of the previous chapter combines an imaging system, a model organism, and a methodology for data analysis. This chapter will describe these in detail. The ideal imaging system must balance many constraints. It must be high resolution so that individuals are well resolved in as large as possible a field. In addition, it must be fast enough to allow one to track and to maintain the identity of multiple individuals in the same arena. However, it must also be inexpensive enough to permit multiple systems to be constructed and run in parallel. Our model organism must be motile, have relatively short generations times, and be have appropriate size and swimming speed to be tracked. Our goal was to image and track microbes swimming at speeds of  $mm/s$  at a resolution of microns for durations of hours across areas of  $mm^2$ . Experiments will necessarily be 10's of hours long to capture multiple generations, and during that time, we cannot allow individual organisms to be hidden or to escape, but should in general be freely moving. During this time, we required that the environment be homogenous, in terms of chemical compositions and physical parameters such as temperature. We also require that the individual identity of each



organism is maintained, even when multiple organism are present.

We strove to minimize costs so that we could construct parallel systems to run many experiments simultaneously and collect population level statistics. For this, we constructed seven replicate digital microscopes based on inexpensive commercially available cameras, fabricated custom chambers in which to isolate and image organisms, and developed custom software to analyze the images and generate trajectories which maintain identity for at least three generations (up to eight individuals). This chapter will describe the technical aspects and performance of our microscopes, introduce important aspects of our chosen model organism, and describe the image analysis and tracking algorithms.

## **2.1 Hardware**

### **2.1.1 Imaging System**

Consumer demand has driven the innovation of high-resolution ( $10^7$  pixels), high-frame-rate (15 Hz) digital cameras. We have used these as low cost alternatives to more expensive scientific cameras in the construction of replicate microscopes for tracking. The current pixel size of the sensors in these cameras is 2-5  $\mu\text{m}$ , which makes them suitable for low to no magnification imaging of microbes. In theory, imaging could be done without any optics [7], but due to limitations in the construction of the imaging sensors we chose, this was determined to be infeasible. The imaging apparatus consists of a light source, a condenser, a sample stage, a focusing relay

lens, and an imaging sensor (Figure 2.1). An image of the sample plane is focused onto the sensor using a 30mm focal length 1x relay lens (Edmund Optics, Barrington NJ). Images are digitized using a consumer webcam (Logitech USA) with a 1600x1200 pixel CMOS sensor that acquires images at 15 fps. The pixel size on the sensor is  $3.3\mu\text{m}$ . The effective pixel size of the image is  $4.25 \pm 0.04\mu\text{m}$ , giving a magnification of 0.78 (i.e. the image is reduced by this factor). Illumination is provided by a single soft white LED (Philips Lumileds, San Jose CA) driven by a constant current of 0.35 A provided by a current regulated driver (LEDdynamics, Randolph VT). Illumination is focused on the sample using a bi-convex lens (Thorlabs, Newton NJ). Movies are compressed and stored to a disk to be processed at a later time.

### **2.1.2 PDMS Chambers**

Each organism is isolated and imaged in a custom fabricated chamber, in a circular arena, 5 mm in diameter and  $230\mu\text{m}$  in depth (Figure 2.2(a)). Chambers are made in-house using soft lithography [37] and fabricated in poly-dimethyl siloxane (PDMS), an optically transparent elastomer (Ellsworth Adhesives, Germantown WI) (See Appendix F). We developed a monolayer valve system, taking advantage of the elastic properties of PDMS, that allowed the isolation of a single individual. This is important because we wanted to isolate a single individual and its progeny for many hours without the possibility of an individual hiding or escaping. Physically, a chamber is created by an 5 mm diameter annulus of PDMS that extends from the ceiling of larger chamber (Figure 2.2(b) upper panel). When this larger chamber is

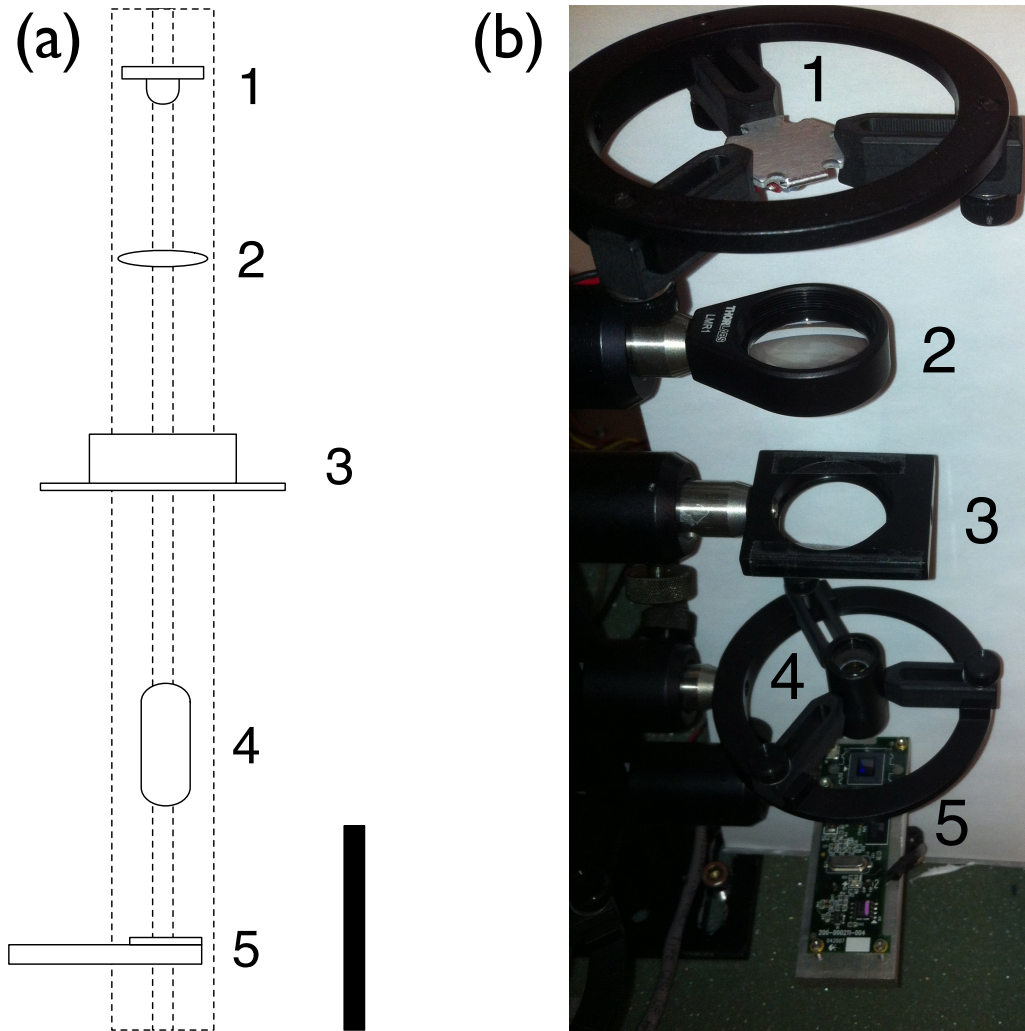


Figure 2.1: **Schematic and Photograph of Imaging System** (a) Schematic of Imaging System, scale bar indicates 50 mm. (b) Photograph of Imaging System, showing (1) Illuminating LED (2) Condensor Lens (3) Sample Stage (4) Relay Lens (5) Webcam Sensor

pressurized, it distends, which lifts this ring from the glass substrate and (Figure 2.2(b) lower panel) allowing organisms and media to flow underneath. Details of the lithography procedure can be found in Appendix B. These PDMS chambers are gas permeable, allowing the free diffusion of oxygen into the system. This can also lead to evaporation, however, if media is maintained in the outer annulus, evaporation in the inner chamber is reduced, and the net evaporation rate is estimated to be  $0.2 \text{ mm}^3$  per day, or about 6% per day.

Because it is important to understand whether behavioral variability arises naturally or in response to changing conditions, we wanted to both ensure that the environments were as homogenous as possible as well as catalogue the physical effects which might constrain the behavior. Because of the quasi-two-dimensional geometry, we first sought to determine an appropriate chamber depth. If the chambers were too shallow, organisms could be physically pinned, or experience significant wall drag, however, if they were too deep, motion in the third dimension would be lost in the projection. To assay this, we conducted experiments in a variety of chamber depths, starting at a shallow depth and increasing it until we saw no evidence of the increase in chamber depth in our measurements of the behavior. This resulted in our choosing  $230 \text{ }\mu\text{m}$  (See Appendix A).

In addition to the interactions with the floor and the ceiling, we also investigated the effect of proximity to the chamber walls. Looking at radial distribution functions of individuals, as well as radial averages of motion parameters, we determined that the chamber could be divided into two regions, termed in the “bulk” and near the “walls”. We found that this to be within  $42.5 \text{ }\mu\text{m}$  of the wall (Appendix A). Distributions of

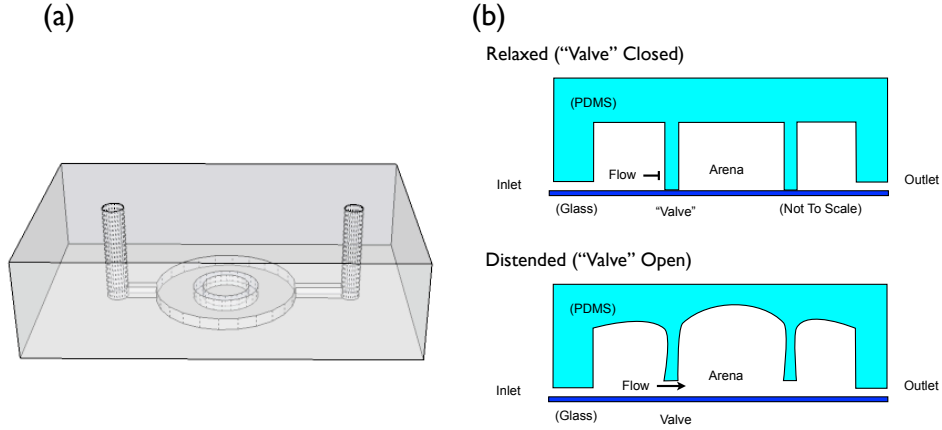


Figure 2.2: **Schematic of PDMS Microfluidic Chamber** - showing a three-dimensional representation of the chamber and (b) showing in cross section, representation of the closed (upper panel) and open (lower panel) configurations.

residence times near the walls show a log-normal distribution with an average of mean interaction of 1 second (Figure A.5). Distributions of the fraction of each individuals lifetime that is spent in the bulk *versus* near the wall show, in general, about 10-20% of the lifetime is spent near the wall (Figure A.6).

Next we sought to determine if there were macroscopic heterogeneities in the chambers. For example, a persistent, uneven distribution of nutrients might cause organisms to favor certain areas of each chamber more than others. To assay this, we looked at spatial cross-correlation functions between individuals in the same chamber and comparing them to individuals in different chambers. We found no significant difference in these correlations indicating no macroscopic heterogeneity that were apparent as a difference in spatial bias.

### 2.1.3 Temperature Control

Temperature is another environmental parameter that might influence behavior, so it was important to control this and ensure that it too was homogenous over the course of the experiment. To achieve this, we designed a temperature control apparatus. A custom amplifier with a linearized thermistor (Omega, Stamford CT) provides an input voltage proportional to the temperature (See Appendix Figure B.1) which is digitized via a LabJack U3 USB DAQ (Labjack, Lakewood CO) interface with a precision of 1.2 mV (4.8 mK). This measured voltage is controlled with a proportional integral control feedback loop using Matlab. The feedback voltage is applied to a amplifier which drives a Peltier element (See Appendix Figure B.2). This PI feedback stabilizes the measured voltage to the set-point with standard deviation of 4.6 mK from the set point across all experiments. The thermometer is calibrated to 50 mK absolute accuracy.

## 2.2 Wetware

### 2.2.1 *Tetrahymena* sp Biology

*Tetrahymena* are unicellular eukaryotic protozoa, approximately 50  $\mu\text{m}$  in length and 15  $\mu\text{m}$  in diameter. Most of the work presented here was done in *Tetrahymena thermophila*. *T. thermophila* are common in fresh water ponds and streams in North America, and have been found as far west as Minnesota and as far south as Florida [12]. Although they are known to be bacterivores, much about their natural lifestyles,

including their preferred prey, and even how they survive the winter, remains a mystery.

*T. thermophila* has been a fruitful model organism for cell biology over the last 50 years. Early work on the cell cycle took advantage of the ease with which cultures could be synchronized to determine causal relationships among different events [15]. Ribozymes, catalytic RNA molecules, were first discovered in *T. thermophila* [28]. Later work in *T. thermophila* was integral to the discovery of telomeres, the caps of DNA that protect its ends from degradation associated with replication [47], as well as telomerase, the enzyme which rebuilds and maintains telomeres [8]. Most recently, the role of histone acetyltransferase and histone acetylation was discovered using *T. thermophila* as a model organism [10].

*Tetrahymena* swim consistently at speeds to up to 1 mm/s by means of many rows of beating cilia. Metachronal coordination of cilia, which refers to the constant phase difference between adjacent cilia, is the means of propulsion and is thought to be coordinated by passive hydrodynamic interactions between the cilia [18]. Cilia are patterned in an average of 18 meridians that extend longitudinally and have a chiral twist [33]. This chirality is preserved during division, and has been used to demonstrate templated cortical patterning, an early example of non-genetic inheritance [35]. *T. thermophila* has been shown to be chemotactic, responding both positively and negatively to a variety of peptide and protein signals [30]. When starved, *T. thermophila* are known to undergo a phenotypic change which results in faster swimming. This change involves the elongation of the body and the growth of a long, caudal cilium [34] and is thought to be a dispersal morph.

*T. thermophila*, like most ciliated protozoa, exhibit nuclear dimorphism. Each organism contains two nuclei, a micronucleus and a macronucleus, essentially a differentiation between germ line and soma as in multicellular organisms. The macronucleus, where all active transcription takes place, is highly polyploid, consisting of an average of 45 copies of each of its 200-300 “autonomously replicating pieces” (ARPs), essentially short chromosomes that generated by fragmentation and rearrangements of the 5 micronuclear chromosomes after conjugation. The micronucleus is diploid and divides mitotically during asexual division. The macronucleus, however, divides amitotically after replication of each of the ARPs, and alternative copies of alleles are segregated at random. The sexual phase of the *T. thermophila* life cycle and consists of the exchange of micronuclear genetic material without cell division between cells that must be different in one of seven different mating types. Mating type is determined genetically at the *mat* locus. Some strains, including *T. thermophila*, can reproduce indefinitely asexually.

*Tetrahymena sp.* have many characteristics that make them attractive as a model system. They are easily maintained in laboratory culture and have a relatively short generation time (four hours). *Tetrahymena sp.* cultures can be maintained long-term in soybean cultures or frozen (See Appendix F). In addition, they swim constantly and are large enough to be imaged easily in our system.

For the majority of experiments, a derivative of *T. Thermophila* Strain CU428 was used, and we refer to this as wild type in what follows. In addition, we performed behavioral analysis on two other strains of *T. thermophila*, natural isolates from New Hampshire and Pennsylvania obtained from the Cornell Stock Center, as well as



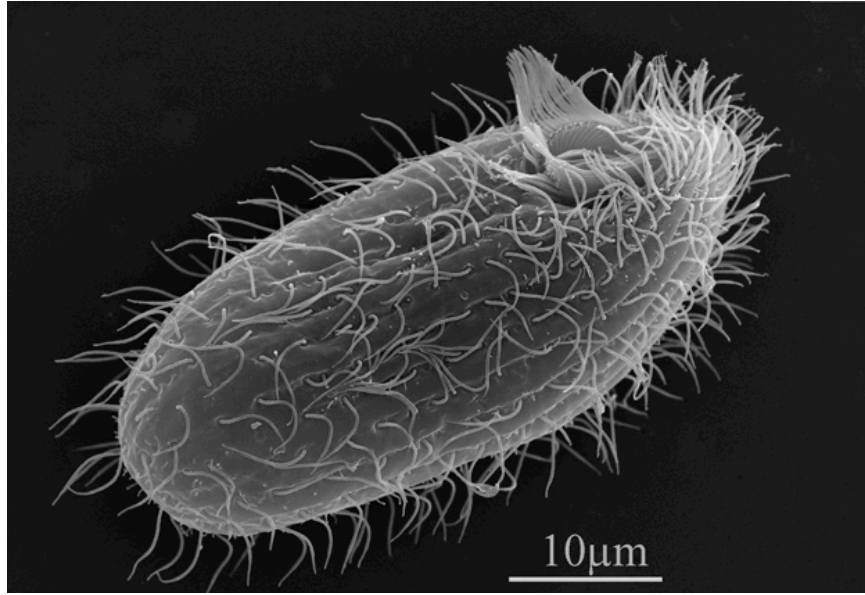


Figure 2.3: **Scanning Electron Micrograph of *T. thermophila*** - taken by Aswati Subramanian, scale bar indicates 10  $\mu\text{m}$ .

another species, *Tetrahymena borealis* (See Appendix F).

### 2.2.2 Growth Media

The base growth medium is SPP medium, 2% proteose peptone, 0.1% yeast extract and 0.2% glucose. To assay the effects of different media on the swimming behavior, we developed a panel of chemical and physical perturbations. These were designed to span a spectrum, from those that mimicked natural environments, to those that were completely artificial. In addition media were chosen that would keep the growth rate as close as possible to the growth rate in standard 1xR media. Cultures were grown at room temperature without shaking for 48 hours for each experiment in the appropriate medium (Table 2.1).

2xR is the is SPP media with twice the concentration of each ingredient (4% pro-

Table 2.1: Abbreviations for Growth Media

1x SPP	<b>1xR</b>
2x SPP	<b>2xR</b>
1xSPP Sterile Filtered + Beads	<b>1xB</b>
Bacterized	<b>Bac</b>
Chemically Defined Media	<b>CDM</b>

teose peptone, 0.2% yeast extract and 0.4% glucose.) [2]. 1xB is the 1xR media that has been filtered with a 0.2  $\mu\text{m}$  sterile filter flask, and then supplemented with 1.57  $\mu\text{m}$  diameter poly-methyl methacrylate beads at a density of 400  $\mu\text{g}/\text{ml}$  of particles (6800/ $\text{mm}^3$ ) [38]. Bacterized media was prepared by growing *Escherichia coli DH5 $\alpha$*  in 1xR media to an OD of 0.6 before sterilizing via autoclave. Chemically defined media was taken directly from [46].

Based on estimates of nutrient concentrations in the media, and nutrient uptakes rates of the *T. thermophila* [2], and the relative volume of the cell to the chamber ( $10^{-6}$ ) we concluded that nutrient depletion would not be significant over the course of our experiment. This is supported by the observation that the carrying capacity of the chamber is many thousands of cells.

## 2.3 Software

### 2.3.1 Image Analysis

At a sampling rate of 15 Hz for an average lifetime of 270 minutes, the movie that results from each experiment is on the order of 250,000 images. Each image is captured

using the webcam software and imported into Matlab (The MathWorks, Natick MA) using the VideoIO toolbox (Gerald Daley). The resulting images are dynamically background subtracted and segmented with a global threshold. Detected objects are filtered based on size and an expected number of objects. The centroid location, as well as the area, eccentricity and orientation of each object is then recorded and stored. Detection is based on the difference between each image  $I(t)$  and the maximum projection  $B(t_0)$  of  $I(t_0)$  and  $I(t_0 + \tau)$ . A new background image is recalculated every  $\tau$  frames. A detection matrix  $M(t)$  is calculated as  $I(t) - B$  and an "object" is defined as a set of connected pixels in  $M$  with an intensity greater the global threshold  $T$ . In this manner, artifacts due to long-time-scale changes in the image over the course of the experiment can be avoided, such as changes in illumination or shifts of the imaged volume, however objects that stop moving completely for more than  $\tau$  frames will not be detected. The parameter  $\tau$  can be adjusted to be longer than the longest period of inactivity observed. For *T. thermophila* we have set this delay to 3000 frames (200 s) and the threshold  $T$  to 0.2 of the maximum pixel value in a background subtracted image. Once images are segmented with this threshold, the centroid, area, orientation, and eccentricity of each connected component is determined in Matlab and stored for subsequent tracking.

### 2.3.2 Tracking

Tracking refers to the association of objects in one frame with the corresponding time displaced object in the next frame. From the segmented images trajectories

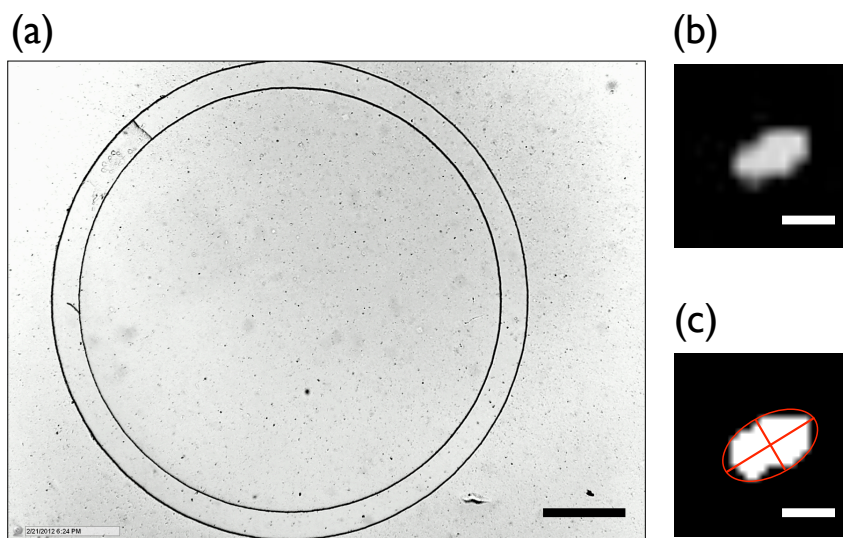


Figure 2.4: **Example Image from Digital Microscope** (a), scale bar (black) indicates 1 mm. (b) shows an expanded view of a background-subtracted image centered on an organism. (c) Shows the resulting thresholded image with the bounding ellipse. Scale bars (white) indicate  $25\ \mu\text{m}$ .

are created using a method derived from [25]. The segmentation algorithm returns a list of objects in each frame. We create a cost function for assigning objects to the same trajectory between neighboring frames. This cost function is based on the Euclidean distance between objects in neighboring frames and an empirically determined weighted difference in object area (See Appendix B). Empirically determined costs for making no assignment for a specific object are also included, this allows the system to be robust to transient occlusions. Objects are first linked pair-wise frame by frame into segments using this cost function, mathematically, this is known as a linear assignment problem. The optimal assignment, the assignment which minimizes the total cost for linking objects from one frame to that in the next, is determined using the hungarian algorithm, a combinatorial optimization algorithm which solves a linear assignment problem in polynomial time. The resulting segments are aggregated and linked in a second linear assignment matching using a cost function which incorporates Euclidian distance and change in median segment area, but also includes the time gap and the change in velocity between the end of one segment and the beginning of another. Because *Tetrahymena* swim at speeds up to 1 mm/s and we image at 15 frames per second, at most, cells are separated by about one body length from frame to frame, this aids in maintaining identity during tracking, even so, each trajectory formed from the joined segments is then checked by hand to ensure individual identity is maintained (See Appendix D).

## Chapter 3

# Analysis of Behavioral Patterns in *T. thermophila*

This chapter describes the analysis of behavioral patterns in populations of *T. thermophila* and related species. These populations were grown in environments which varied chemically and physically. Trajectories were recorded as time-series of spatial locations in two dimensions and transformed to corresponding time-series of linear and angular speeds. These time-series were discretized and two-dimensional histograms corresponding to the frequencies of the discretized values were used as representations of the underlying trajectories and the similarity between two behaviors is determined using a relative entropy based metric between such histograms.

### 3.1 Preliminary Observations

Generation times for individuals in populations grown in different environmental conditions varied by as much as 35% of the mean across all populations (Figure 3.1). The number of frames that comprises each trajectory gives the generation time of that individual, where divisions are inferred from the splitting of blobs in subsequent frames

of the movie. Interestingly, while doubling the concentrations of nutrients in the standard growth medium affects cell size (data not shown), it does not reduce the generation time, in fact the average generation time is slightly longer in this condition. Shorter generation times can be achieved by the bacterization process. The chemical composition of 1xR and 1xB is identical, thus a physical change generates a significant increase in generation time.

Generation time correlations were seen between the two sister cells which are the product of a biological division. Furthermore, generation time and post division size are perfectly correlated (Figure 3.2) with the smaller sister always having a longer subsequent generation time than the larger sister. However, while there are a number of studies which detail asymmetric division in *T. thermophila* [29], determining the mechanism for this correlation is beyond the scope of this work.

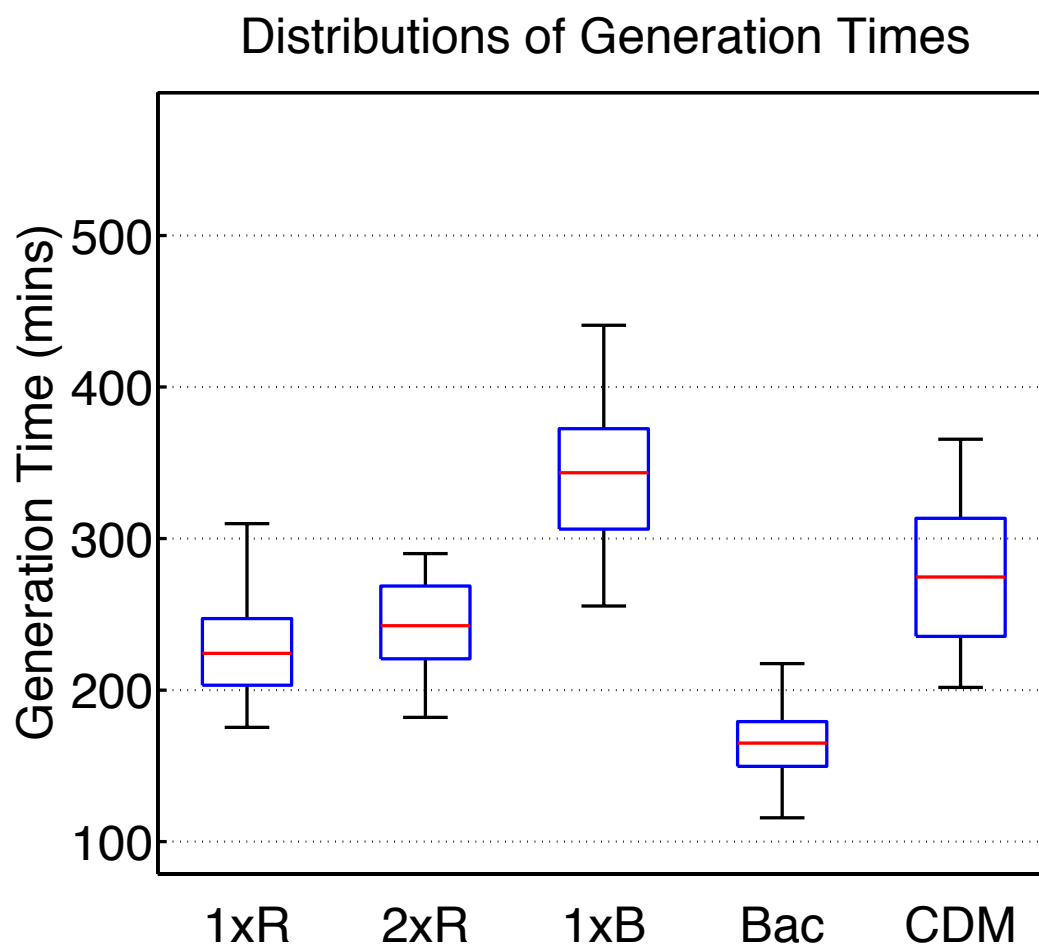


Figure 3.1: **Distributions of Generation Times in 5 Environments** - for  $N=30$  individuals in 5 different environments. Red lines are medians, bars show quantiles and whiskers show extrema.



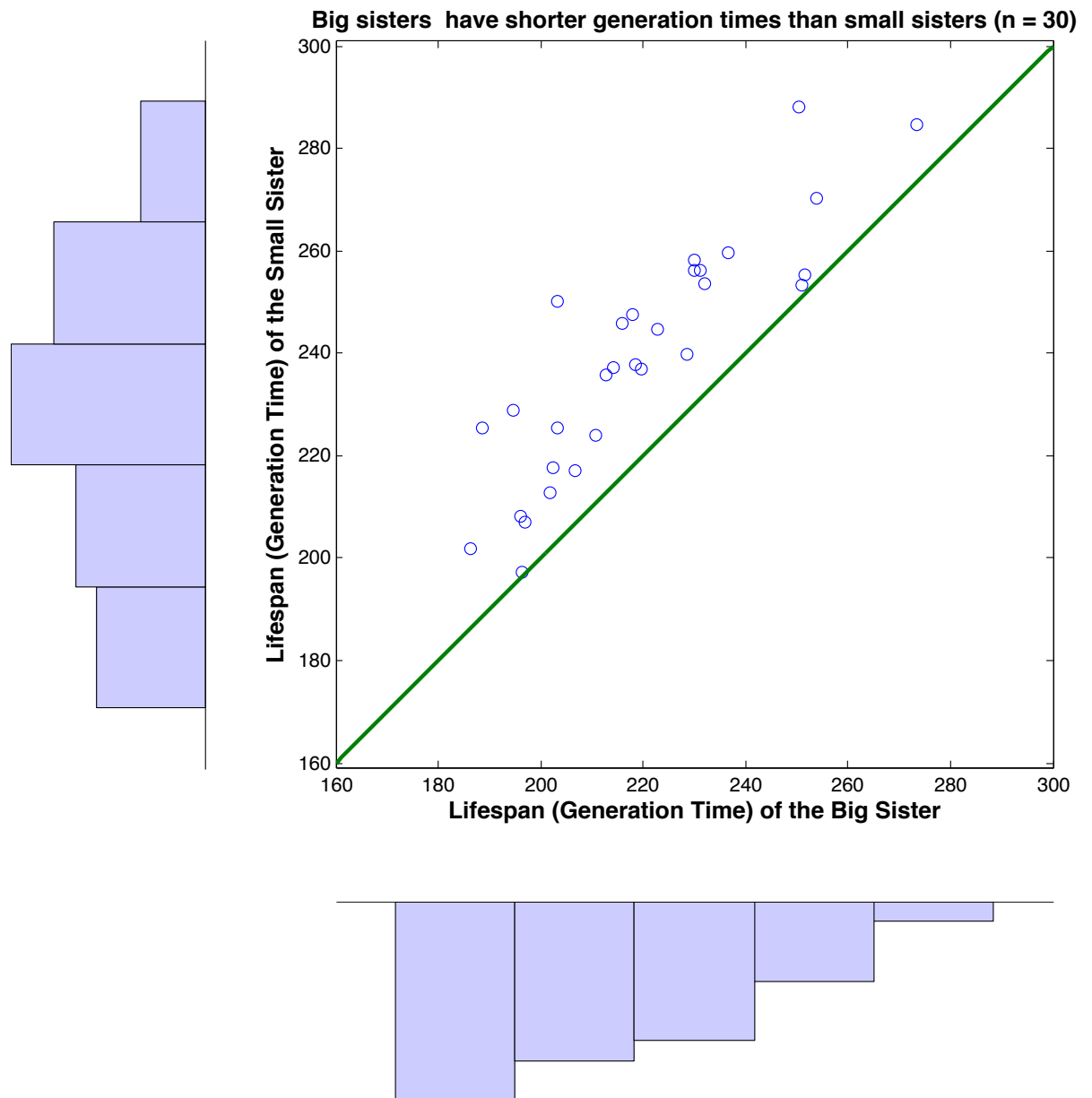


Figure 3.2: **Generation Times of Sisters are Perfectly Correlated with Post Division Size** - Data is shown for individuals grown in 1xR medium.

## 3.2 Transformation of Data

Trajectories consist of (x,y) coordinates as a function of time at a frequency of 15 Hz (Figure 3.3). For the analysis presented here we omit regions of the trajectory where the individual is within 42.5  $\mu m$  of the chamber boundary (see Appendix A). Let  $v = (v_x, v_y) = ((\Delta x/\Delta t), (\Delta y/\Delta t))$ . linear speed is then the norm of v, and angular speed is given by

$$sign(v(s) \times v(t)) * \cos^{-1} \frac{(v(s) * v(t))}{(|v(s)| |v(t)|)}$$

Where  $s = t + \Delta t$ , and  $\Delta t = 1$  frame or 1/15 s. Transforming the data from spatial locations to component velocities relies on assumptions of spatial isotropy and homogeneity, which we have confirmed (Appendix A).

## 3.3 Quantitative Representation of Behaviors

To answer the questions posed in the introduction relies on being able to directly measure differences between behaviors. Behavioral differences have traditionally been classified by scoring, which breaks sequences of actions into pre-recognized classes called stereotypes. However, scoring does not permit comparisons between stereotyped classes and breaks down when the boundaries between classes become unclear.

Here we introduce statistical methods to measure the similarity of behaviors that are described as distributions of their underlying actions, avoiding scoring altogether. In our methodology, an observable action is speed and turning angle pair. Behaviors at different time scales are thus represented by sequences in the time series of the

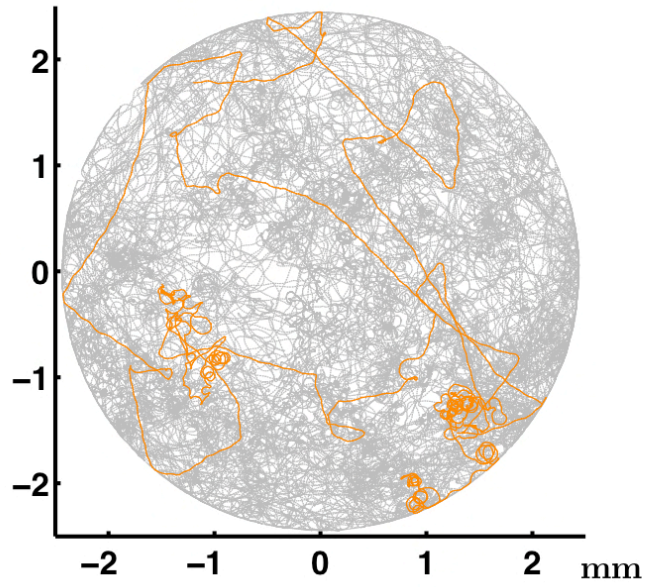


Figure 3.3: **Example Trajectory for a Single Individual** - grey trace shows the trajectory of a single individual over its entire lifetime ( $\approx 240$  min), while the orange trace shows a 1.1 minute segment. Cursory examination indicates at least two modes of motion, a faster ballistic type, and a slow diffusive type.

appropriate duration. Thus the behavior on the longest time scale for an individual is the entire lifetime sequence of speed and omega. In addition, we can look at arbitrarily shorter segments of the trajectory. As described above, the average trajectory consists of on the order of 250,000 points, described in both speed and omega. Because data becomes sparse in such high dimensional space, we first sought to provide a meaningful reduced representation of the data. For this we chose joint distributions of the frequencies of linear and angular speed pairs  $P(\omega, |v|)$ . The use of two-dimensional distributions retains important aspects of speed and turning angle correlations. We used empirical histograms to estimate these distributions. In this work behaviors are defined as histograms drawn from non-overlapping 1.1 minutes intervals from the time series. We have chosen the length of intervals so that these two-dimensional histograms, approximating the distributions  $P(\omega, |v|)$ , are well populated. Histograms quantify frequencies of observed actions and provide a standardized output across behaviors (for a more detailed explanation, see Appendix C).

### 3.3.1 Measuring Statistical Distance

There are many statistical metrics available to compare distributions. One of the simplest is the Kolmogorov-Smirnov statistic, however this is only applicable to one-dimensional distributions. When we looked at results obtained using it on one dimensional distributions of speed alone, we found that many interesting features of the turning angle correlations were hidden. One of the simplest distances we can imagine is the L1 distance between the histograms, given by  $\delta_{L1}(p, q) = \frac{1}{2} \sum |p(A) - q(A)|$

Surprisingly, this is equivalent to twice what is called the total variation distance. The total variation distance is given by  $\delta(p, q) = \sup_{A \in \Omega} |p(A) - q(A)|$ . Intuitively, this metric determines the action whose frequency is most different between the two histograms and defines the distance between them as this maximum. Because the L1 distance on probability measures is mathematically equivalent the total variation distance, it is dominated by the differences in a single action, and the relative differences between all other actions are ignored. This manifests as a failure to detect differences using the L1 metric that can be detected using other metrics.

We could have used a variety of statistical distance measures, many of which fall into a class called f-divergences, which includes the total variation distance, and other measures such as the Kullback-Leibler divergence and the Hellinger distance. However, many of these do not satisfy the requirements of a true metric, so were excluded from consideration. Furthermore, many statistical distances do not have clear interpretations. Ultimately, we chose to use a metric called the Jensen-Shannon divergence, because it has a clear interpretation, captures variation that other measures do not, and satisfies the properties of a metric. This quantity provides an answer to our original questions about how representative are single measurements. The Jensen-Shannon divergence answers the following question: given multiple observations, if we choose an action randomly from one of the observations, how much information does the identity of that action give about the observation from which it came. Intuitively, if two observations have few actions in common, it is very likely that we can determine from which observation an action was selected. Formally, this can be stated as follows; given that we observed an action  $s$ , how much information

do we gain about whether  $s$  came from observation  $P$  or observation  $Q$ . If we have chosen  $s$  at random and take the average of the information gain over all choices of  $s$ , this information gain can be computed and is called the Jensen-Shannon divergence [31], [17].

The Jensen-Shannon divergence satisfies the requirements of a metric, one of which is symmetry. Symmetry is important because the similarity between two behaviors should not depend on the order in which they are listed. If  $p_i$  and  $q_i$  are the frequencies of action  $i$  for two behaviors, then the JS divergence between those behaviors is given by the following equation:

$$D(p, q) = \frac{1}{2} \left[ \sum_i p_i \log \left( \frac{p_i}{m_i} \right) + \sum_i q_i \log \left( \frac{q_i}{m_i} \right) \right]$$

Where  $m_i = \frac{p_i + q_i}{2}$ . With this metric, the distance between two behaviors that have no actions in common is one and the distance between two behaviors that have the same actions, identical in frequency, is zero. Behaviors with more actions in common are less distant.

### 3.4 Patterns of Individual Behavioral Change

Individuals show behaviors that change during their lifetimes, and these changes can be adaptive. It is clear that such changes are not mediated by mutation and selection. While some of the changes are initiated by changing environments, there is a great deal of change that occurs in homogeneous environments. In the following sections,

we will present our efforts to quantify and catalogue the types of behavioral changes that can occur in individuals in homogeneous environments, and to measure the time scales over which these changes occur.

### 3.4.1 Changeability

Swimming behaviors, represented by distributions  $P(\omega, |v|)$ , can be compared to one another by computing the Jensen-Shannon divergence,  $D(\bullet|\bullet)$ , between their histograms. We can compute such a distance for every pair of distributions at different times  $t$  and  $t'$ . The set of all such distances, which we call the changeability, is represented by a matrix because of the discretization of  $t$  and  $t'$  and is given by the following equation, where the histogram that represents the behavior of individual  $N$  at time  $t$  is denoted  $P^N(t)$ .

$$C^N(t, t') = D(P^N(t)|P^N(t'))$$

(In what follows we will drop the symbol  $D$ , denoting Jensen-Shannon divergence, for convenience). The process of computing a changeability matrix is diagramed in Figure 3.4(a), and the resulting matrix for a single individual is shown in Figure 3.4(a) upper panel.

Changeability matrices allow us to measure temporal trends of behavioral change within an individual's lifetime. Using these matrices, we present data that catalogues the diversity of such trends in a variety of chemical environments for different strains of *T. thermophila*. Our data indicate that these patterns themselves are highly variable.

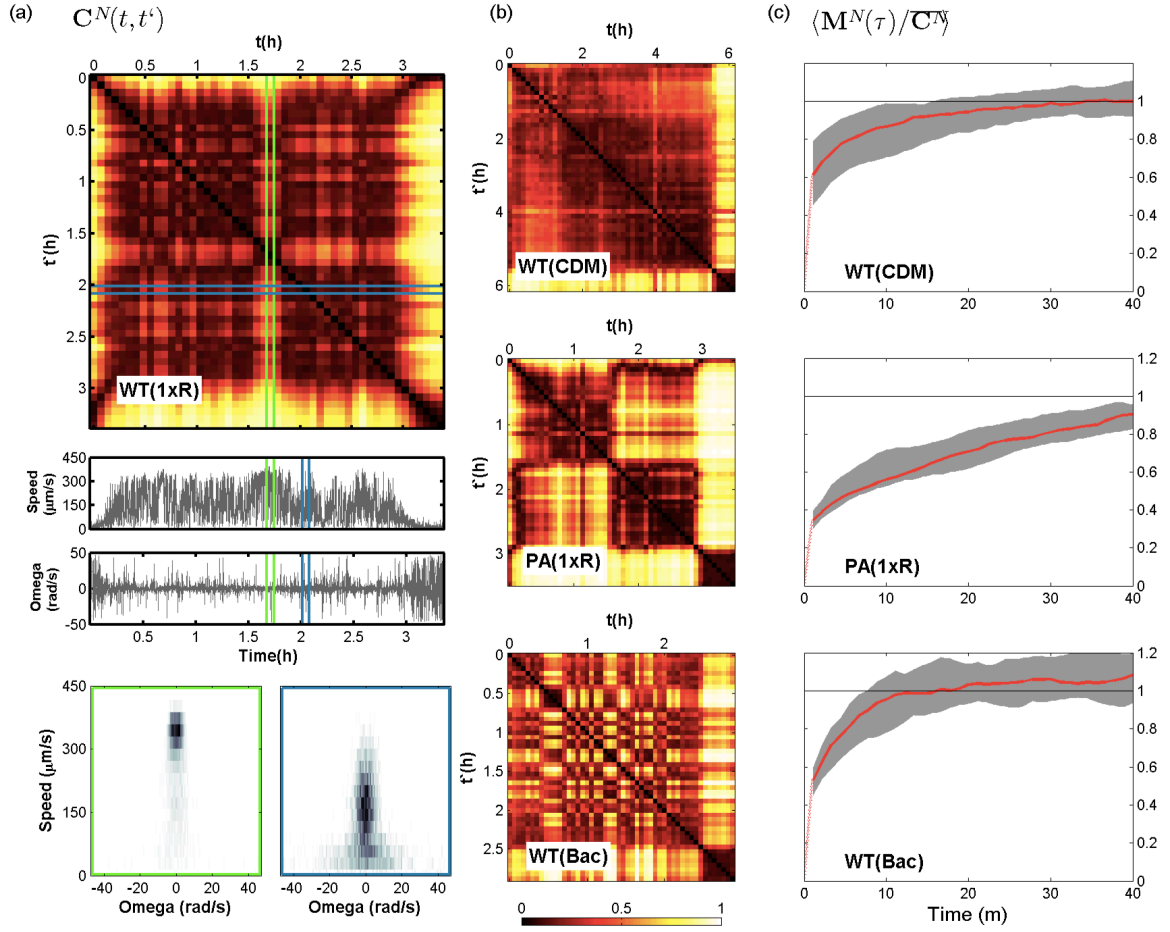


Figure 3.4: **Changeability Matrices Show Patterns of Individual Behavioral Change** - (a) Indicates schematically how a histogram (lower panel) are each associated with a single window (denoted by green or blue outlines) and correspond to a given row and column of the changeability matrix. Middle panels show the speed and omega data used in constructing the changeability matrix in the upper panel, the lower axis of all three plots are identical. White labels in the lower right corner denote strain identity and environmental condition. (b) Shows some examples of different characteristic changeability profiles 3 different illustrative conditions. The dimensions of each matrix are scaled to each individuals lifetime so they do not correspond panel to panel. (c) Shows the population averaged normalized memory  $\langle (M^n(\tau))C^n \rangle$  for each of populations ( $n=30$ ) shown in (b). Red lines indicate the median normalized memory (they go trivially to the origin) and gray areas indicate the 0.1 and 0.9 quantiles. Peri-division windows are discarded for this memory calculation. Characteristic time scales can be estimated as the time scale at which the median normalized memory reaches 1.0.



For example, we observe that behavioral changes can accrue continuously over an individual's lifetime (Figure 3.4(b) upper panel).. This type of behavioral dynamics indicates that differences between behaviors more distant in time are larger and that this trend persists throughout the lifetime of the individual. This constant accumulation of behavioral difference in time is akin to the notion of aging and can be regarded as an example of behavioral aging [23].

In contrast, large behavioral changes can appear suddenly (Figure 3.4(b) middle panel). This individual clearly shows a dramatic behavioral transition that occurs at about the midpoint of its lifetime. This abrupt change in the behavior of an individual is reminiscent of life-cycle stages [36], [24], although relating this observation to internal events in *T. thermophila* life cycle lies beyond the scope of this work. The second sudden transition that occurs near the end of the lifetimes can be associated with cellular division.

While the first two examples show structured patterns of behavioral change, we have also observed individuals that exhibit large behavioral changes that occur frequently and at seemingly random intervals throughout the whole lifetime (Figure 3.4(b) lower panel).

### 3.4.2 Behavioral Memory

The various transitions shown in Figure 3.4(b) are examples from a wide spectrum of observed dynamics. In order to summarize the varied dynamics we have observed, we present a quantity called the behavioral memory. The behavioral memory can be

associated with the following question: given that a behavior is observed at some time  $t$ , how different a behavior, on average, will be observed in that individual after some time  $\tau$ ? This quantity is a form of an autocorrelation function of behaviors. Memory quantifies the speed of behavioral change and is given by the following equation:

$$M(\tau) = \langle \overline{(P^N(t)|P^N(t+\tau))} \rangle$$

The over bar represents an average over all times,  $t$ , for an individual and the angle brackets denote the average value of  $M^N(\tau)$  for all individuals in a population.

The memory reflects our statements about how rapidly behavioral patterns change and this quantity is shown in Figure 3.4(c) for the populations from which the individuals in Figure 3.4(b) were drawn. The slowly changing individual in Figure 3.4(b) middle panel, characterized by long periods with little change, was drawn from a population with a behavioral memory of 1 hour or 30% of the lifetime, while the rapidly changing individual in 3.4(b) lower panel was drawn from a population with a behavioral memory of  $\approx 10$  minutes which is 5% of those individuals lifetime. The individual in Figure 3.4(b) upper panel comes from a population with  $\approx 40$  minute behavioral memory. While this is comparable to the slow dynamics in absolute time, relative to the average lifespan, it is only 10%.

The examples of changeability matrices in Figure 3.4(b) were chosen to highlight familiar behavioral patterns, e.g. progressive aging, stages of life, etc. In general, our data show that the temporal structure of behavioral change varies depending on the individual. This has important implications for the design on experiments that seek

to characterize behaviors.

First, if large behavioral changes accrue continuously over time, any sub-lifetime measurement will not be representative of the full lifetime behavioral repertoire (Figure 3.4(b) upper panel). Further, large and abrupt transitions in behavior indicate that even if one measures behavior for an hour and finds low changeability, a claim that such a measurement is representative of that individual's behavior would be erroneous (Figure 3.4(b) middle panel). Therefore, full lifetimes must be measured to fully characterize the behavior of an organism.

Furthermore, our results show that the behavioral memory, which quantifies the average time-scales of behavioral change among many individuals, is dependent on both the environment and on genetic makeup. Thus even if one measures full-lifetimes and finds short memory in one environment or for one genotype, short measurements cannot be confidently made in a different environment or in a mutant. Therefore, full lifetime measurements must be replicated for each environmental or genetic perturbation.

### **3.5 Behavioral Variability in Populations**

It is known that a single genotype can exhibit a variety of phenotypes. Survival or extinction depends on the adaptation of such phenotypes to the environment. Selection acts on this distribution of phenotypes, as survival of the individual is secondary to survival of the reproducing population. Thus it is important to understand how phenotypes like behavior can vary among individuals in a population. In the fol-

lowing, we will introduce methods that use our quantitative behavioral descriptions to look at the variability of behaviors between individuals in populations. We will show that behaviors can change when organisms are grown in different environmental conditions, an example of what is classically called phenotypic plasticity. In addition, our data show that behaviors can, in some conditions, persist from one generation to the next. We present a measure of this persistence which gives a simple measure of heritability which represents the similarity of behaviors as a function of the number of generations. In contrast to population genetics, we do not impose an underlying model of genetic relatedness.

### **3.5.1 Individuality**

In addition to computing changeability, we can use our formalism to compare the behavior of one individual to that of others, a measure which we call individuality. When comparing individuals, it can be simpler to compare their average of their behaviors, rather than the each of behaviors separately. While it is unclear how to average stereotyped behaviors, our approach allows averaged behaviors to be represented as the average of their respective histograms. Thus, the single histogram generated from the sequence of all observed actions represents an individuals average behavior and the JS divergence between two such histograms is the distance between those individuals. With such histograms, we can compare the average behaviors between individuals and look for differences resulting from genetic or environmental changes.

With measurements and characterizations of full lifetime behaviors, we can extend our analysis to quantify the differences in such behaviors between individuals. Histograms (Figure 3.5(a) lower panel) can be generated from an individuals entire lifetime sequence of speed and turning angles. These histograms approximate the distributions,  $\overline{P^N(\omega, v)}$ , that describe an individuals average lifetime behavior. In general, the measurement of the differences between the average lifetime behavior ( $\overline{P^N}$ ) of individual N and another individual M, can be represented as a matrix and is given by what we call the individuality matrix (Figure 3.5(a) right panel):

$$I(N, M) = (\overline{P^N} | \overline{P^M})$$

Just as changeability quantifies the behavioral repertoire of a single organism the individuality matrix  $I(N, M)$  quantifies the extent of the behavioral repertoire of individuals in a population. While the changeability matrix is naturally ordered by time, the individuality matrix can be ordered arbitrarily, for example, in the matrices we present we have chosen to order individuals descendent from a single progenitor (a family) together and further to group families grown in the same environment together (Figure 3.5(a)).

Individuality matrices allow us to observe patterns of behavior in populations. We observe that behavioral differences between individuals can reflect differences in the environments or in the relatedness of those individuals.

### 3.5.2 Plasticity

Phenotypes that differ as a result of environmental differences are said to be plastic. We have found conditions for which behavior in *T. thermophila* is strongly plastic (Figure 3.5(b) middle panel). Individuals grown in bacterized media (Bac) show behaviors that are distant from all behaviors seen in individuals grown in standard media (1xR). However, it is clear that even in these strongly plastic conditions, the individuality distributions are overlapping. Between other sets conditions, however, individuals can show very similar behaviors. In cases where there is no plastic response to an environmental change, behavioral differences are characterized by largely overlapping distributions (Figure 3.5(b) right panel).

### 3.5.3 Heritability

Similarly, the correlation of phenotypes with relatedness is known as heritability. Heritable phenotypes show smaller differences between more closely related individuals and we observe conditions in which behavior is heritable (Figure 3.5(c) middle panel). In CDM, two individuals drawn from the same family tend to exhibit similar behaviors, but these behaviors are distinct from those exhibited by another family. Again, note that these differences are characterized by distributions of distances, and that these distributions are overlapping, even when behavior is strongly heritable. Families with largely different behaviors can still have individuals with similar behaviors.

Heritability is similar to behavioral memory, but it extends the idea to include changes over multiple generations, thus strong heritability is akin to a slow de-

correlation of behaviors with each generation (Figure 3.5(c) middle panel). However, in different conditions, behaviors can de-correlate quickly, even in a single generation (Figure 3.5(c) right panel).

These observations of behavior in populations have important implications for experiments. First, to characterize the behavioral differences between populations, it is essential to measure the behavior of enough individuals so that individuality distributions are statistically well characterized. This is because, regardless of plasticity or heritability, differences between populations are characterized by distributions and these distributions can be overlapping. Therefore, with only a few measurements, differences may be uncharacteristically large or small since they can be affected by statistical fluctuations, e.g. come from the tails of individuality distributions.

In addition, one must be careful to avoid measuring correlated populations of individuals. We show that behavioral differences between individuals can reflect the relatedness of those individuals. Thus, a population derived from a single progenitor or a single clutch of eggs may show very different properties than one derived from a group of isogenic founders or a collection of eggs from a group of isogenic mothers. Furthermore, our data show that if such correlations exist, it is not clear, a priori, for how many generations they will persist. This correlation time must be measured and incorporated into experiments that rely on the assumption of uncorrelated populations.

Figure 3.5: **Behavioral Variability in Populations** - (a) This panel demonstrates how two individuals (indexed as 2 and 5, green and blue outlines) in this condition (1xR) contributes to the individuality matrix (boxed entry). Strong plasticity (b, middle panel), differences between subgroups to be individuals grown in different environmental conditions, is shown as histograms of individuality divided into within and between condition subgroups. Blue and green bars are drawn from the within condition individuality ( $I_{(e=e')}$ ) matrix (inset, blue and green triangles). The black bars show the histogram of the individuality between individuals in different conditions, which are drawn from the between condition portion of the individuality matrix (inset, black square). Strong heritability (c, middle panel), differences in individuality between more and less related individuals, is shown as histograms between related (black bars) and unrelated individuals (green bars). The corresponding parts of the individuality matrix are shown with black and green (inset). Data is organized in a way that puts related individuals on the block diagonal. (b and c right panels) Indicate conditions of weak plasticity and weak heritability, observed as overlapping distributions for within and between group comparisons.



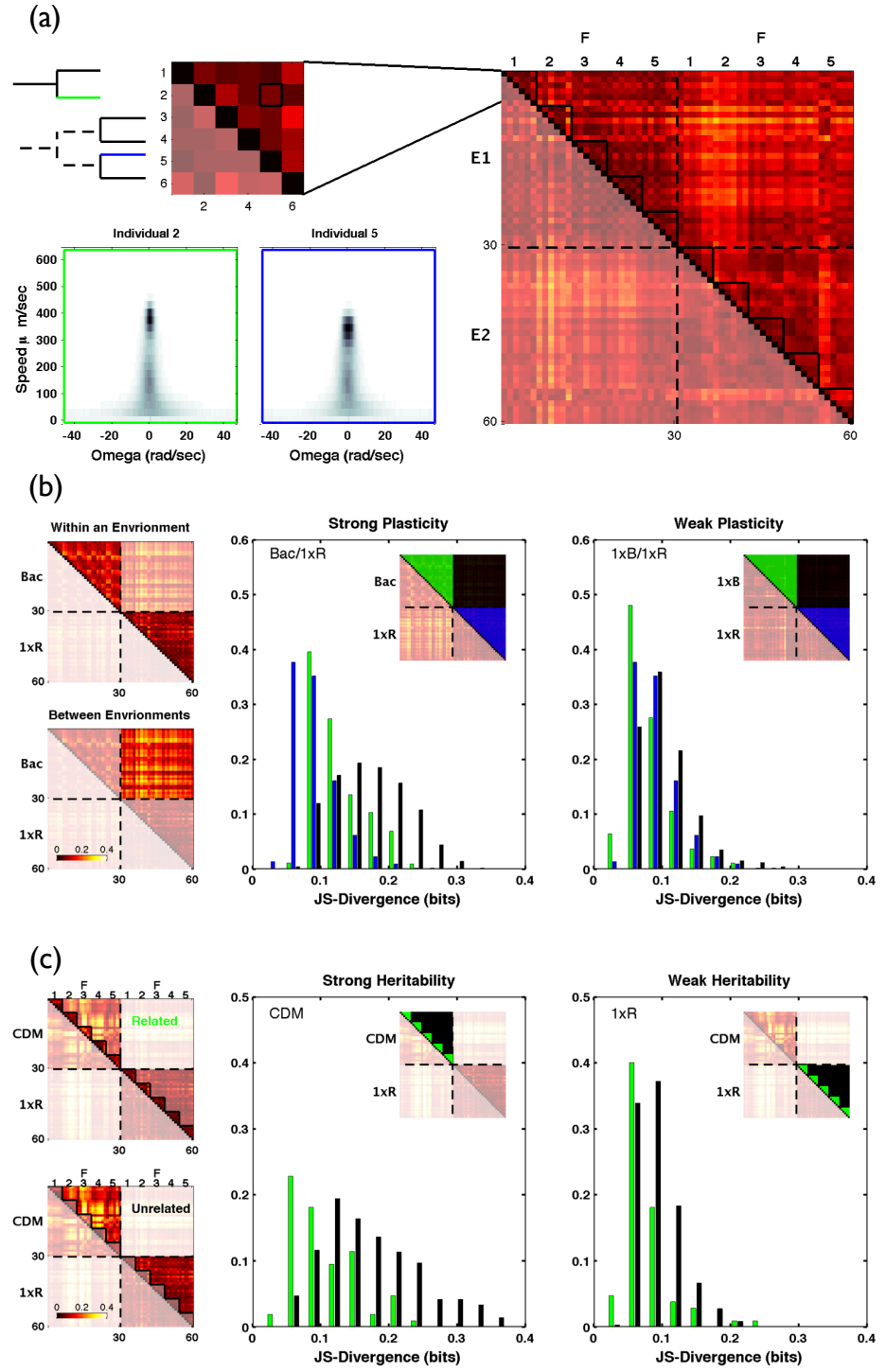


Figure 3.5: Behavioral Variability in Populations

Table 3.1: Quantities which describe behavioral dynamics

Operation	Notation
Time Average	$\bar{\bullet}$
Population Average	$\langle \bullet \rangle$
JS Divergence	$(\bullet   \bullet)$
Quantity	Equation
Changeability	$C^N(t, t') = (P^N(t)   P^N(t'))$
Memory	$M^N(\tau) = \overline{(P^N(t)   P^N(t + \tau))}$
Individuality	$I(N, M) = (\overline{P^N}   \overline{P^M})$
Plasticity	$1 - (I(N, M)_{e=e'} / I(N, M)_{e \neq e'})$
Heritability	$1 - (I(N, M)_{f=g} / I(N, M)_{f \neq g})$

### 3.6 Dimensionality Reduction

Dimensionality reduction refers to the transformation of high-dimensional data (inputs) into a meaningful representation in a reduced dimension (outputs). By meaningful we mean, in general, that nearby inputs are mapped to nearby outputs and distant inputs are mapped to distant outputs. Often, the reduced dimensional representation corresponds to the intrinsic dimensionality of the data. Thus, such a procedure can reveal the minimum number of parameters needed to reconstruct the observed variation in the data. This is appealing for a system where the ‘primitives’ are unknown, such as in our behavioral system. If we can find a meaningful reduced dimensional representation of the data, this may help to uncover the simple behavioral primitives from which complex behavioral patterns are generated. Ideally, we would then be able to find biological features, for example, genetic elements or coherent network states, that correspond to these primitives. The following section will give a brief overview of techniques for dimensionality reduction, as well as present our

attempts to represent sequences of behaviors in reduced dimensional space, which we will call behavioral space.

### 3.6.1 Multidimensional Scaling

Multidimensional scaling (MDS) refers to a collection of techniques that attempt to find a distance preserving map from high-dimensional inputs to reduced dimensional outputs. Multidimensional scaling methods attempt to find this mapping by minimizing the difference between distances in the embedding, and the actual matrix of pairwise distances, called the distance or dissimilarity matrix. Classical linear dimensionality reduction techniques, based on vectorial representations of the data, generates a linear map between the input and output space. If multidimensional scaling is performed using the euclidean distance between objects in the input space as the distance matrix, MDS is equivalent to eigenvalue decomposition. If the distance matrix is generated by computing correlations between the input vectors, MDS will yield the same results as principle components analysis (PCA). MDS can be extended to generate nonlinear mappings by changing the way the distance matrix is calculated. For example, if we compute only the  $k$  nearest neighbor distances between input vectors using a euclidean metric, and then determine longer distances using by finding the shortest path, this results in a non-linear technique known as Isomap [48]. Our choice of distance metric, the Jensen Shannon divergence also relaxes the assumption of a linear input space. Defining distances between behaviors using the JS divergence necessitates the use of MDS as a dimensionality reduction technique. If we wanted to use linear methods such as PCA, we would first have to find a vectorial representation of the data. Two candidates for such a representation might be the histograms themselves, or a vector of preselected features, which could be chosen based on intuition or based on a model, such as a gaussian mixture or hidden Markov model. These choices would result in the requirement that all observed histograms be either linear

combinations of characteristic "eigen"-histograms, or a linear combinations of model parameters. We found that this is not a reasonable assumption. Multidimensional scaling allows the generation of embeddings directly from the similarities computed by JS divergence, i.e the changeability and individuality, that we have presented.

Dimensionality reduction problems can often be formulated as an optimization problem, and for MDS, the outputs  $\phi_i \in \mathbf{R}^m$  are chosen to minimize the normalized stress function given the inputs  $x_i \in \mathbf{R}^d$  with  $m < d$ :

$$\varepsilon = \frac{\sum_{ij} (\|x_i - x_j\| - \|\phi_i - \phi_j\|)^2}{\sum_{ij} \|x_i - x_j\|^2}$$

This can be achieved by singular value decomposition of a pairwise distance matrix. Intuitively, this procedure attempts to generate a configuration of points in a given  $m$  dimensional space, such that the euclidean distance between those points is as close as possible to that defined in the distance matrix. For our purposes, that distance matrix will be a changeability matrix, a set of changeability matrices, or an individuality matrix. Thus, the distance between behaviors in the embedding should faithfully represent the Jensen Shannon divergence between the histograms that represent those behaviors.

### 3.6.2 Evaluating Embeddings

In principle an embedding in the intrinsic dimension of the data should have a stress of zero, any non-zero stress is an indication of insufficient dimensionality. In practice, uncertainty and under sampling can lead to non-zero stress even when at the intrinsic dimensionality of the data. We can evaluate how much residual stress is relieved as we increase the dimensionality by plotting the residual stress as a function of the dimensions in the embedding space, known as a Scree plot (Figure 3.6). At least for single strains in a variety of conditions, two dimensional embeddings seem to

be sufficient, although when considering more than one strain or species together, a third dimension may be appropriate. For visualization purposes, embeddings will be presented in 2 dimensions in what follows.

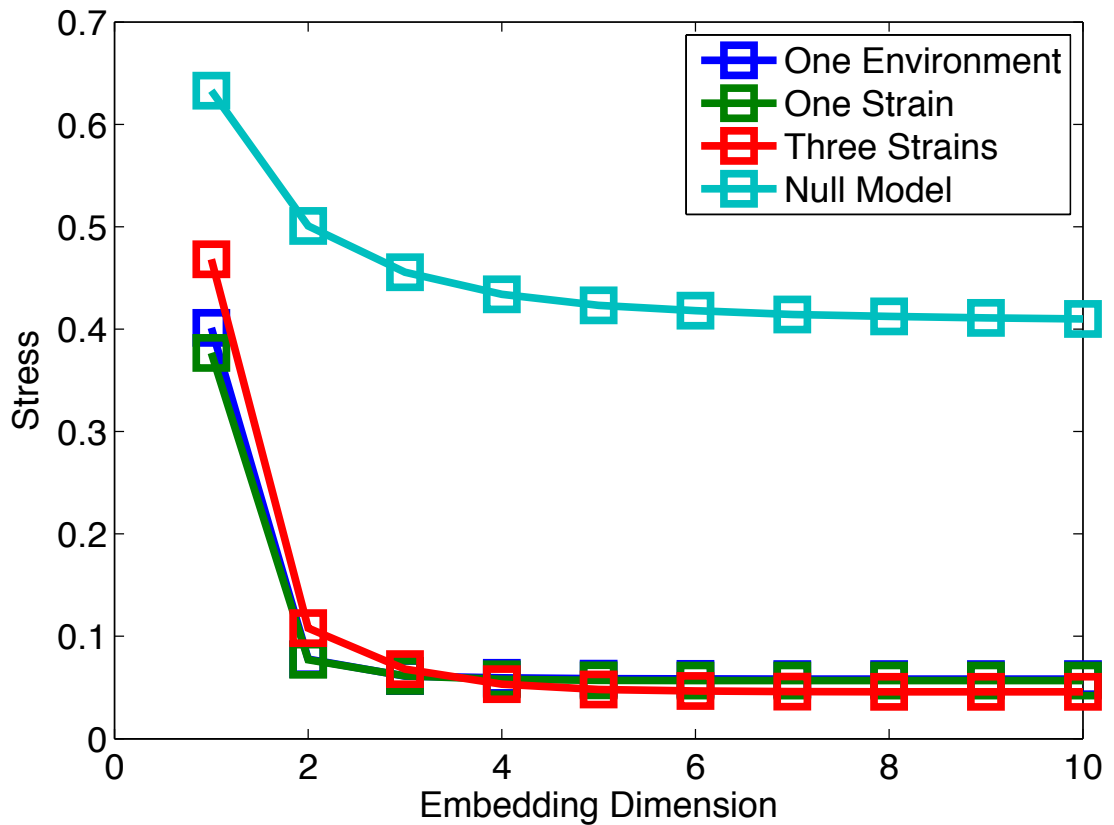


Figure 3.6: **Scree Plot of Embeddings for Different Populations:** Residual stress as a function of dimension for 30 WT individuals in one Condition, (blue), 150 WT individuals in 5 different environments (green), and 90 Individuals from 2 different strains and one different species in 1xR (red). A null model (teal) was generated by randomizing the entries in the changeability matrices of individuals.

### 3.6.3 Trajectories in Behavioral Space

Properties of changeability and behavioral memory are captured in behavioral space trajectories. Returning to the three changeability matrices from the examples in Figure 3.4 Features of changeability matrices are captured by these dynamic trajectories. In each case, the division state (red points) is clearly separated from rest of the trajectory, reflecting the large differences between  $P(\omega, |v|)$  histograms during division and the rest of the life. In addition, the two distant behavioral states of the individual in 3.7(b) are clearly seen as distinct clusterings of points in the trajectory. The concept of behavioral memory is also nicely captured in these behavioral space trajectories. In each, temporally adjacent points are connected by line segments, thus the length of a single line segment in relation to the overall spatial extent of the trajectory is related to the behavioral memory. The individual in 3.7(a), which was shown to have a long behavioral memory, shows individual steps that are short compared to the total extent of the trajectory, while the individual in 3.7(c), which was shown to have short memory, shows more ergodic exploration of its behavioral extent. Behavioral space trajectories, in addition to allowing one to visualize dynamic properties of the individuals, also provide a convenient way of visualizing the effects of perturbations on the behaviors of entire populations. Individuality matrices can be embedded in the same way as changeability matrices, but recall that individuality matrices are generated by comparing average lifetime behaviors. By embedding changeability matrices of entire populations, the similarity of sub-lifetime behaviors between individuals from different populations can be determined. Thus while embeddings of individuality matrices will reveal the similarity of average lifetime behaviors, embedding populations of changeability matrices can show, for example, the existence of some behaviors in a population which are not seen in a different population. The embeddings for entire populations are very dense, so to aid visualization, we will present contour plots generated from applying a kernel density estimator to the original embedding (Fig-

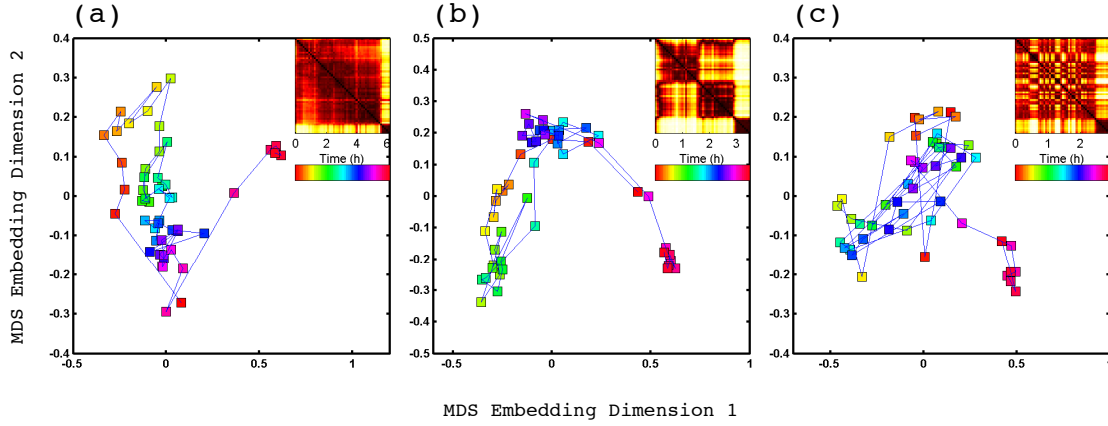


Figure 3.7: **Behavioral Trajectories from Multidimensional Scaling of Changeability Matrices:** Trajectories in behavioral space for the three individuals shown in Figure 3.4. Color of points along each trajectory indicates the time at which that behavior occurred. Insets show the changeability matrices that generated these embeddings.

ure 3.8). Between 1xR and 2xR, the repertoire of behaviors is largely overlapping, however, there are some behaviors in 1xR that are dissimilar from all behaviors in 2xR, and that this dissimilarity is not as large as that with the division state. In the following section, we will present behavioral space representations of environmental and genetic perturbations on populations.

### 3.6.4 Effects of Environmental and Genetic Perturbations

#### Environments

If populations grown in different environmental conditions are embedded in the same space, plasticity (Figure 3.5(b)) is reflected in the overlap of densities of those populations. Conditions of weak plasticity will result in largely overlapping densities, while strongly plastic conditions will show some regions of non-overlapping density (Figure 3.9(a)). The densities that represent 1xR and 1xB are largely overlapping

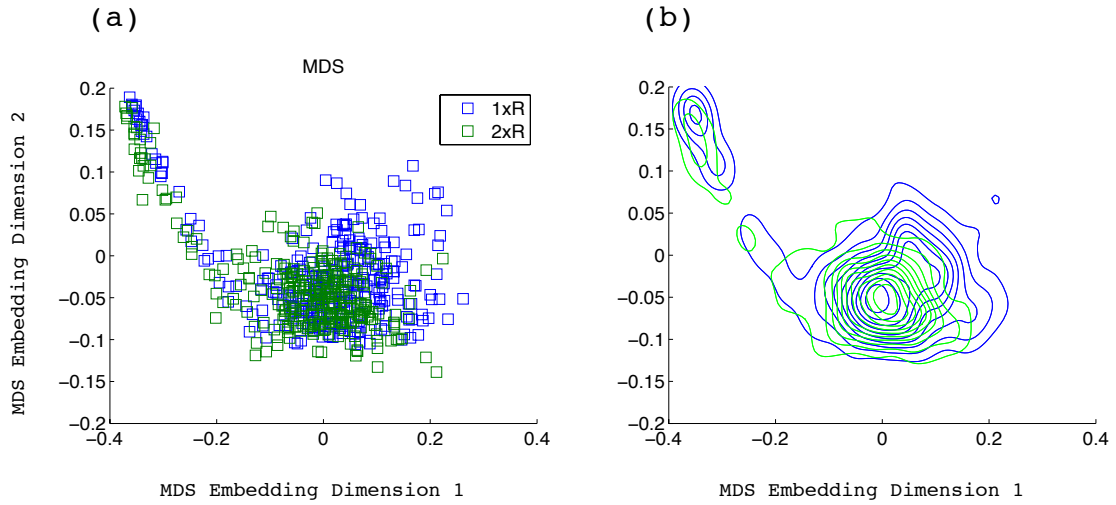


Figure 3.8: **Raw Embedding and Contours of Kernel Density Estimate:** of Populations from 2 Environmental Conditions - (a) shows the raw embedding of populations (n=30) from 2 different environmental conditions, 1xR (blue) and 2xR (green). (b) shows the contour plot of the kernel density estimate.

(Figure 3.9(a) blue and green), reflecting the weak plasticity as determined from the individuality matrix. The densities of 1xR and Bac (Figure 3.9(a) blue and red) show large areas that do not overlap, reflecting the observation of strong plasticity. Thus, in addition to the large differences observed between the average behaviors of individuals grown in 1xR and Bac, it is also apparent that some behaviors exhibited during the lifetimes of individuals in Bac are not seen at all in 1xR.

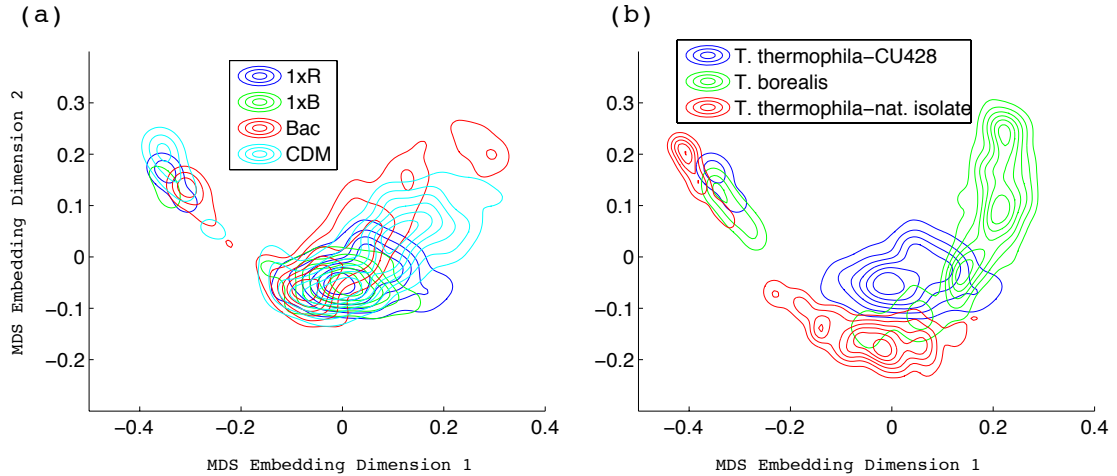
## Genetic Backgrounds

Changes in the genetic background, either behaviors in a different strain of *T. thermophila* or in a different species of *Tetrahymena*, *Tetrahymena borealis* (See Appendix F), are reflected in behavioral space and are generally more dramatic than those induced by environmental perturbations. There is very little overlap between any of these densities, however, as might be expected, the average differences between the



two strains of *T. thermophila* are smaller than those between those of different species.

Although it is beyond the scope of this work to provide a detailed connection between environmental signals, genetic changes, and behavior, these preliminary observations indicate that there may be a meaningful low dimensional representation of behavior, and that it may provide insight into its underlying biological system. In the next section, we will present some preliminary data where artificial selection is applied and the resulting changes are recorded in behavioral space.



**Figure 3.9: Environmental and Genetic Perturbations are Apparent in Behavioral Space Embeddings:** (a) Growth in different environmental conditions can reveal novel behaviors, in general, behavioral repertoires are largely overlapping as compared to (b) behavioral repertoires among different genetic backgrounds grown in the same conditions.

## 3.7 Dimensions in Behavioral Space

### 3.7.1 Perceptual Mapping

With non-parametric dimensionality reduction methods such as generalized MDS, it is often difficult to get an intuition for what the principle axes mean. One way to address this is to generate what is called a perceptual map. Essentially, we superimpose the actual data set (in this case, the two-dimensional histograms) onto the behavioral space projection and look for correlations (Figure 3.10). Based on the perceptual map, we can observe properties which correlate with the large-variance axis (red) and the short variance axis (blue). The long variance axis seems to correlate most strongly with the relative residence time in a high-velocity mode of motion. This agrees with the observation that within an individual, most of the differences are due to state residence time, as the location of the high-speed peak is fairly stationary. The large differences between high-speed peak locations arise when comparing multiple individuals, and indeed, different individuals tend to vary along the blue axis (data not show). Plotting one-dimensional histograms of speed alone for histograms that lie along this axis allows the correlations to be seen more clearly (Figures 3.11 and 3.12). It is clear from these that the histograms vary mostly in the high-speed density. The first example (Figure 3.11) was chosen to highlight the largest changes in state density, however, because the red example does not show any high-state density at all, it has no high-speed mode, and thus we cannot evaluate whether the high-state mode variability is low, as would be expected from the orthogonal arrangement of these histograms to the blue axis in Figure 3.10. A clearer example of the low variability in the high-state mode, despite changes in the state density, is shown in Figure 3.12. In this example there is significant density in the high-state mode, showing that despite changes in the low state density, the location of the high-state peak does not change very much. Lastly, Figure 3.13 shows a clear example of changes in the location of

the high speed peak as we move along the blue axis in Figure 3.10. We can see the variance in the location of the high-speed peak as points move in this dimension in behavioral space is much greater than in the previous cases where points were orthogonal to this axis.

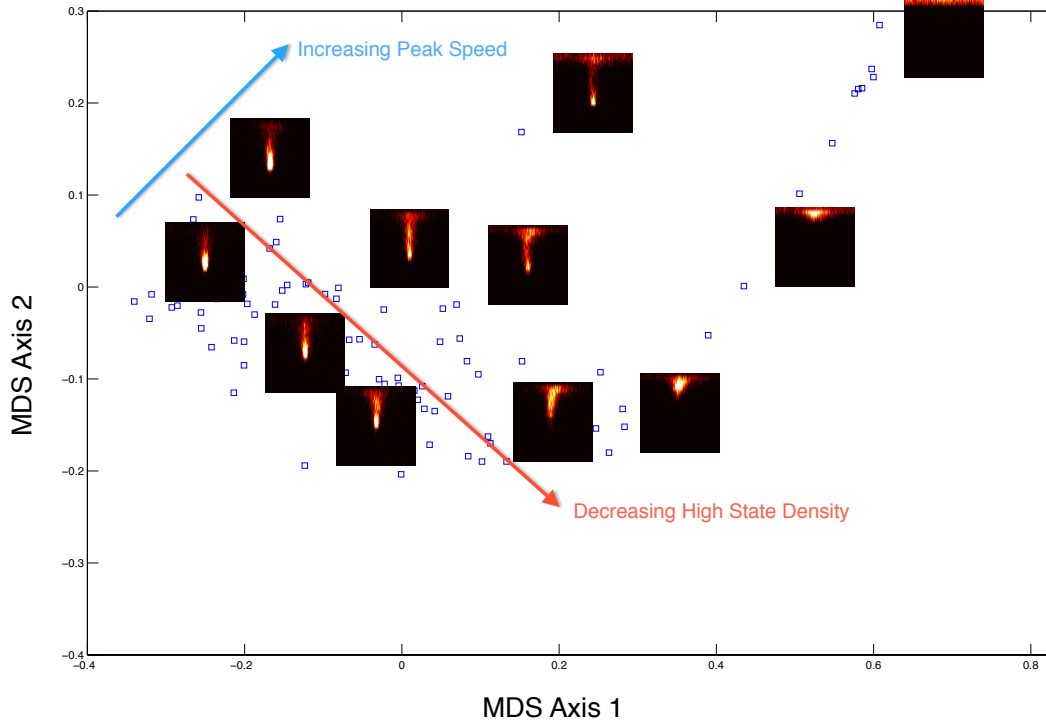


Figure 3.10: **Perceptual map of a single trajectory in 2 dimensions** The tail to the upper right is the division state (high density in the slow state). Subjective assessments of the "meaning" of the two dimensions are indicated by red and blue arrows, they are roughly (blue) an increase in peak speed location, and (red) a change in the relative state density. The change in state density can be seen on these heatmaps, but the peak speed location is harder to visualize

### 3.7.2 Evidence from Selection in Behavioral Space

The previous section shows that one of the directions in the MDS projection may be related to the location of the high-speed mode. To explore this connection further,

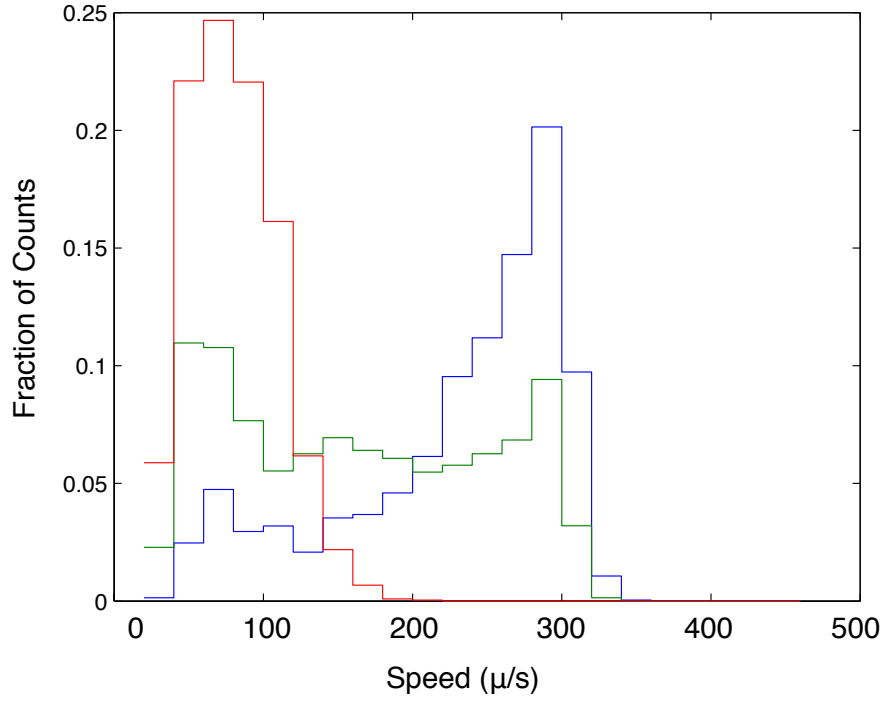


Figure 3.11: **Changes along the red axis correlate with changes in high-speed mode density** An example of changes in state density for histograms that lie above the red axis in Figure 3.10, here it is clearest in the low state density. The red curve comes from  $(-0.2, 0.15)$ , the green curve from  $(0, 0.04)$ , and the blue from  $(0.18, -0.15)$ . These points showed are orthogonal to the peak speed location (blue) axis, and the relatively small variation in the peak speed location can also be observed in the data

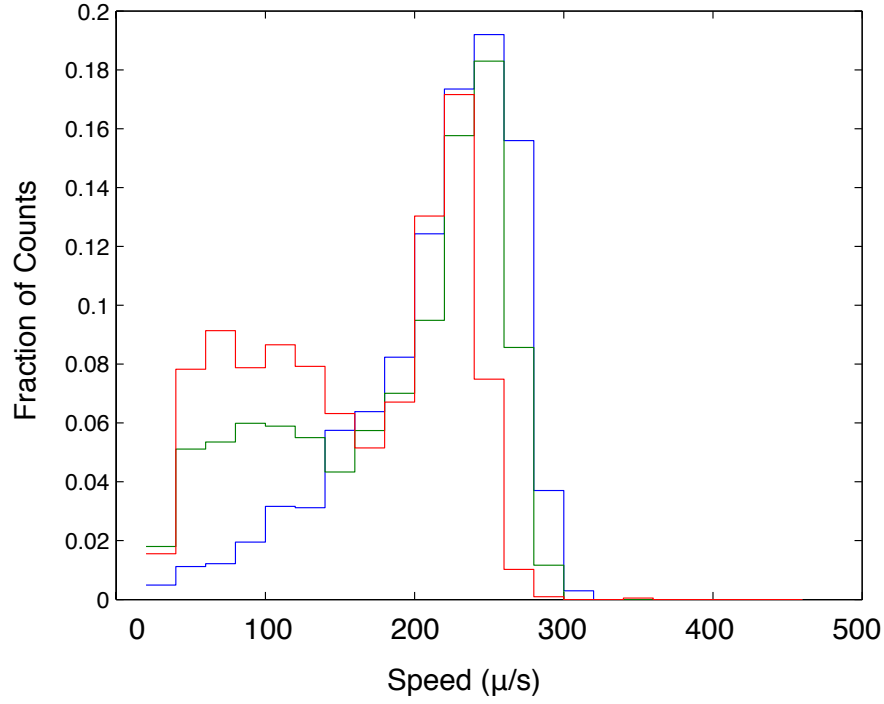


Figure 3.12: **Changes along of the red axis show little change in high-speed mode location** An example of changes in state density for histograms that lie below the red axis in Figure 3.10, again it is clearest in the low state density. The red curve comes from  $(-0.24, 0.05)$ , the green curve from  $(-0.18, -0.05)$ , and the blue from  $(-0.04, -0.16)$ . These points also are orthogonal to the peak speed location (blue) axis, and the relatively small variation in the peak speed location is easily observed

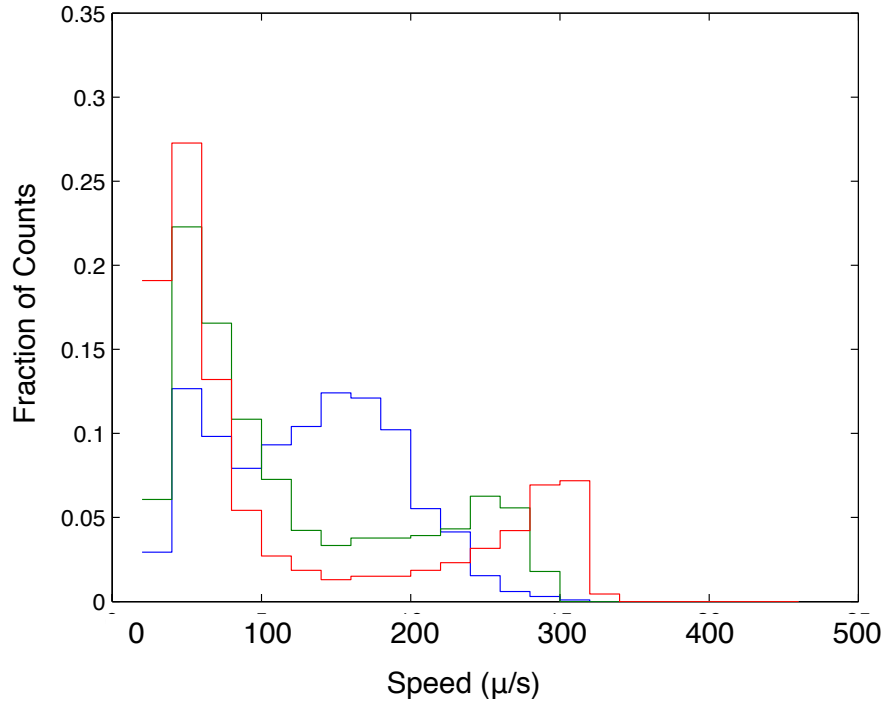


Figure 3.13: **Changes along the blue axis correlate with changes in high-speed mode location** Three histograms drawn from representative points along the blue axis in Figure 3.10, showing changes in the location of the high state mode. The blue histogram lies at approximately  $(-0.2, -0.2)$ , the green from  $(0.1, 0)$  and the red from about  $(0.175, 0.175)$

we took advantage of the asymmetric division to apply artificial selection to a lineage of *T. thermophila*. The relative position of the high speed peak in swimming velocity is perfectly correlated with generation time (Figure 3.14). Thus by removing the longer lived daughter cell after each division, a single lineage that is selective for both larger cells, and for slower peak swimming speed can be selected. To accomplish this, we designed a microfluidic divide with 2 chambers separated by a monolayer valve, from which we could selectively remove the small daughter after each division. This was accomplished by manual removal every 4 hours; long-term experiments of this nature will require a more automated approach. For clarity, the average lifetime behaviors for each individual in CDM, the medium in which the selection experiment was done, were used to generate an individuality matrix, which was then embedded in a two dimensional subspace (Figure 3.15(a) red plus). The individuals in the lineage subjected to selection ((Figure 3.15(a) blue circle), are labelled according to their position in the lineage (1-9). The histograms of only linear speed are shown as a heat map and the decreasing position of the high speed peak is apparent (Figure 3.15(b)). We can see that selection for this parameter drives behavioral changes that move roughly diagonal to the embedding axes and there is a clear correlation between the selective pressure and the resulting behavioral change.

This chapter has presented a set of statistical tools to characterize patterns of behavioral change both exhibited by a single individual through its lifetime, and by groups of individuals in populations. We have shown that we can quantify magnitudes and rates of behavioral change in individuals, using measures we call changeability and behavioral memory, and that these measures reflect familiar notions of individual behavioral variability, such as aging and life-cycle stages. We have also shown that we can use the same framework to compare the behavioral repertoires of individuals in populations. This has allowed us to quantify behavioral differences between populations as a result of environmental differences (plasticity) as well as measure

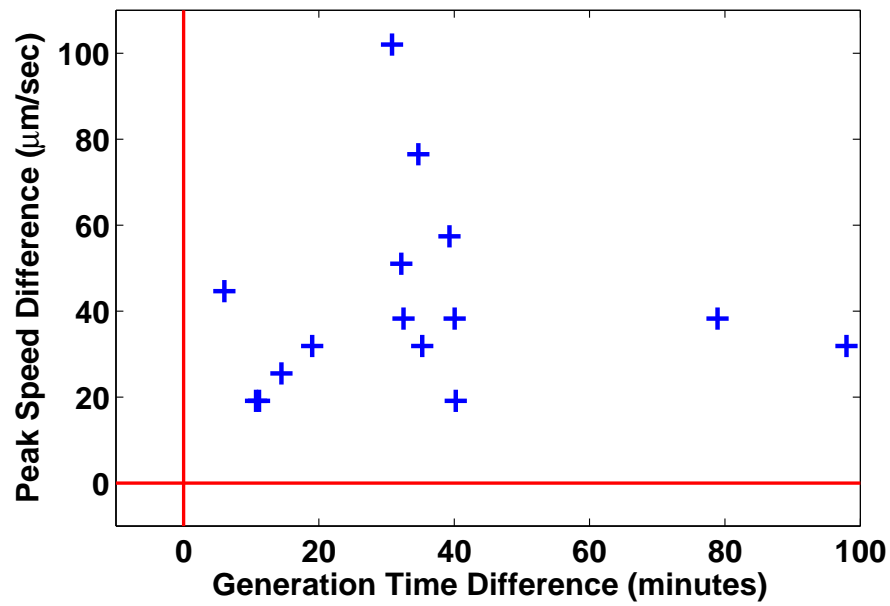


Figure 3.14: **Longer lived sisters after divisions always have a high-speed mode that is faster** This figure shows the 100% correlation ( $n = 30$ ) between generation time difference and peak speed difference among sisters after division. All points have positive values for both the difference between the generation times and the high-speed peak locations.



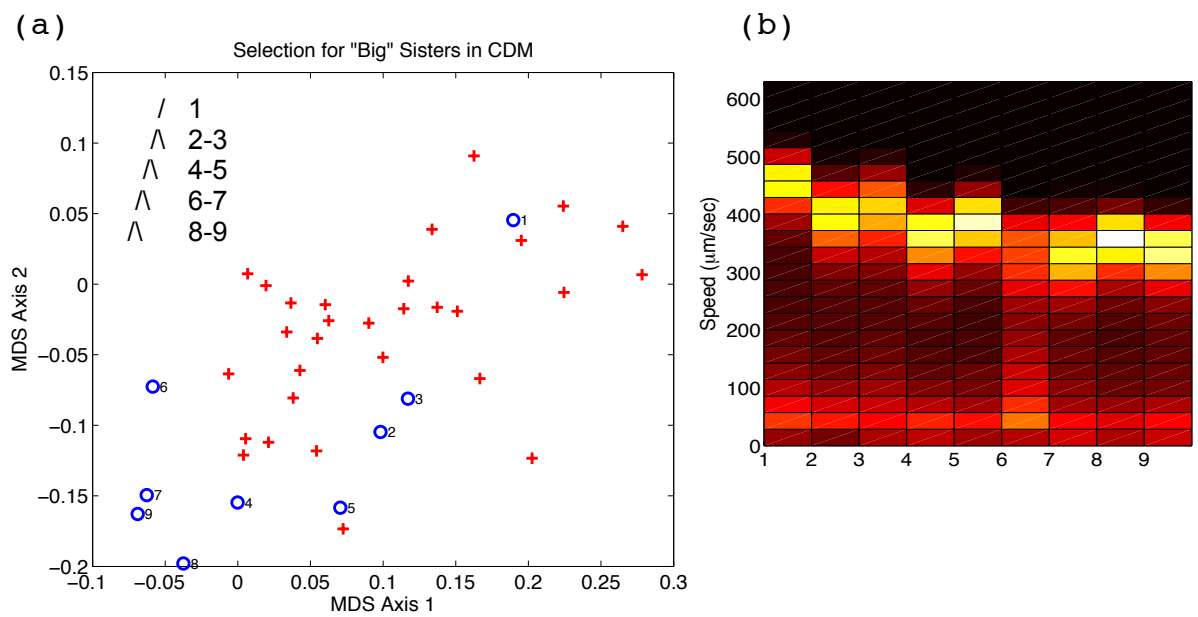


Figure 3.15: **Selection Induced Behavioral Changes can be Observed in Behavioral Space:** (a) Selection for the larger daughter cell after each round of division drives behavior from the upper right (blue circle (1)) to the lower left (blue circle(9)). Because of the size/speed correlation, subsequent generations are selected to have a slower peak swimming speed (b).

decorrelation of behaviors as lineages diverge (heritability). Finally we have presented evidence that despite a system with many degrees of freedom, there may be an appropriate low dimensional representation of the behavior that is biologically relevant.

# Chapter 4

## Conclusions and Future Directions

### 4.1 Conclusions

Microbial swimming is a convenient model system from which to collect high-resolution, long duration spatial data and to characterize motions. We set out to answer some simple questions quantitatively, given that we had trajectories of motions for individuals for their entire lifetimes. Changeability addresses the question of how to measure the variety of behaviors that are exhibited by an individual during its lifetime. Individuality extends that measurement to differences between individuals. Changeability matrices elucidate patterns of behavioral variability in individuals, and behavioral memory quantifies the average time scale of that variation. While behavioral memory measures the persistence of behaviors within a lifetime, heritability measures the behavioral correlations that persist from one generation to the next. We have shown that the variety of behaviors, both in individual and populations, and the persistence of those behaviors, both within and between generations, depend on both environmental and genetic parameters. Although causal relationships have not been established, simply observing these differences has important implications for exper-

imental design. We have also presented some reduced dimensionality representations of swimming behavior, a representation we call behavioral space. Behavioral space has proven to be a convenient way to visualize concepts such as behavioral memory and changeability. Proximity in behavioral space, as measured by euclidean distance, corresponds to similarity of behaviors as measured by JS divergence. Thus similar behaviors will cluster in behavioral space, and clusters that overlap will have many behaviors in common. This allows us to visualize the effects of environmental and genetic perturbations on behavior. We have found that while environmental changes induce some behavioral variability, genetic changes seem to drive even greater diversification. In an attempt to relate directions in behavioral space to relevant biological features, we performed artificial selection for peak swimming speed. Selection for the location of the high speed swimming peak, as performed by size selection, drove behaviors largely along a single direction in behavioral space. Simplifying descriptions in biology, such as the bacterial growth laws, rely on the quantitative measurements of effective parameters, such as growth rate. Because patterns of movement are much higher dimensional in general, any hope of finding simplifying relationships relies on finding appropriate coarse-grained descriptions. Whether behavioral space will provide such a description remains to be investigated.

## 4.2 Future Directions

All the experiments described here were performed using *Tetrahymena sp.* as a model organism. As such, nuclear dimorphism is an obstacle to characterizing or introduc-

ing genetic variation. Future work, which has begun in our lab, will focus on *E. coli* as a model system, where isogenic populations are easy to maintain and precise genetic variants are easily obtained or made. Despite the fact that *E. coli* is 10 times smaller, we have been able to fabricate the appropriately scaled microfluidic devices and incorporate a magnifying objective to image them with the same basic system.

If simplifying descriptions of behavior exist, they will rely on first finding the proper effective parameters. We propose that this will require assaying many environments and many genetic backgrounds, perhaps even many different organisms. Dimensionality reduction might be a fruitful approach to begin to look for such parameters. However, reduced dimensional descriptions might not capture the features of behavior that are biologically relevant. We would like to investigate the relationship between similarity measures, dimensionality reduction techniques, and known biological perturbations, such as targeted genetic changes. In particular, we are interested in whether there are similarity measures that capture more information about the geometry or topology of trajectories.

Taken together, the measures presented here allow for the quantitative study of many previously qualitative ideas. Changeability allows measurement of concepts such as aging, and life cycle stages. Behavioral memory and heritability shows persistence of behaviors over two orders of magnitude, from minutes to days. Measurements of plasticity show explicitly how changing environments drive populations to exhibit different behaviors, and all of these features, even heritability, show environmental dependence. We have presented representations of behavior in two dimensions, which provide a convenient method of visualizing patterns of behavioral change. Prelim-

inary selection experiments have provided evidence that these representations may have biological relevance. If simplifying descriptions of biological systems exist, their discovery will require quantitative measurements, replicate experiments, and systems where internal and external parameters are easily measured or controlled.

# Appendix A

## Chamber Homogeneity and Boundary Effects

Transformation of trajectories from 2D spatial locations to time-series of linear and angular velocity components relies on the assumption of a spatially homogenous and isotropic environment. In addition, the confined 2-dimensional geometry of the chambers used for this experiments required that we investigate the effects of interactions with the boundaries. In addition, because the media in the chamber is not refreshed over the course of the experiment, we sought well to investigate the interactions with any spatial or temporal heterogeneity in the growth media. These sections address the following questions; do interactions with the ‘walls’ of the chamber have an effect on the swimming behavior and if so, how long lasting is it, and are there significant interactions with the upper or lower boundaries of the chamber (‘floor’ and ‘ceiling’), are there spatial heterogeneities in the media that lead some locations in the chamber to be preferred (e.g. food patches), and are there temporal changes in the media that could alter behavior (e.g. nutrient depletion).

## A.1 Chamber Depth

To assess the effects of chamber depth on swimming, 3 sets of experiments ( $N = 18$  individuals) were done at 3 different depths,  $45\mu m$ ,  $85\mu m$  and  $230\mu m$ . We see an effect in peak swimming speed, which is the same for  $85$  and  $230\mu m$  geometries, but decreased for  $45\mu m$  geometry (Figure A.1). In addition, JS divergence between individuals in either  $85$  or  $230\mu m$  and  $45\mu m$  are significantly greater than between  $85$  and  $230\mu m$ . This observation is in accordance with experimental and theoretical estimates of wall drag on swimming ciliates [50], (Figure A.2) This also agrees with our measurement of radius of interaction with the walls.

In addition, if we look at the similarities in the behaviors between  $230\mu m$  and the other chamber depths using the JS divergence metric we have described, we see significant differences between  $230$  and  $45$ , but not between  $230$  and  $80$  (Figure A.3).

## A.2 Walls

Differences in swimming are observed as a function of proximity to the walls of the chamber. Using radial averages of residence time, linear speed, and angular speed, we were able to define a threshold that demarcates the ‘bulk’ from the ‘wall’ at  $42.5\mu m$  (Figure A.4).

Figure A.6 shows the distributions of the fraction of the lifetime spent on the wall which was variable for the five conditions tested. The distribution of the length of wall events was also condition dependent, and the average distribution for each condition is shown in Figure A.5.



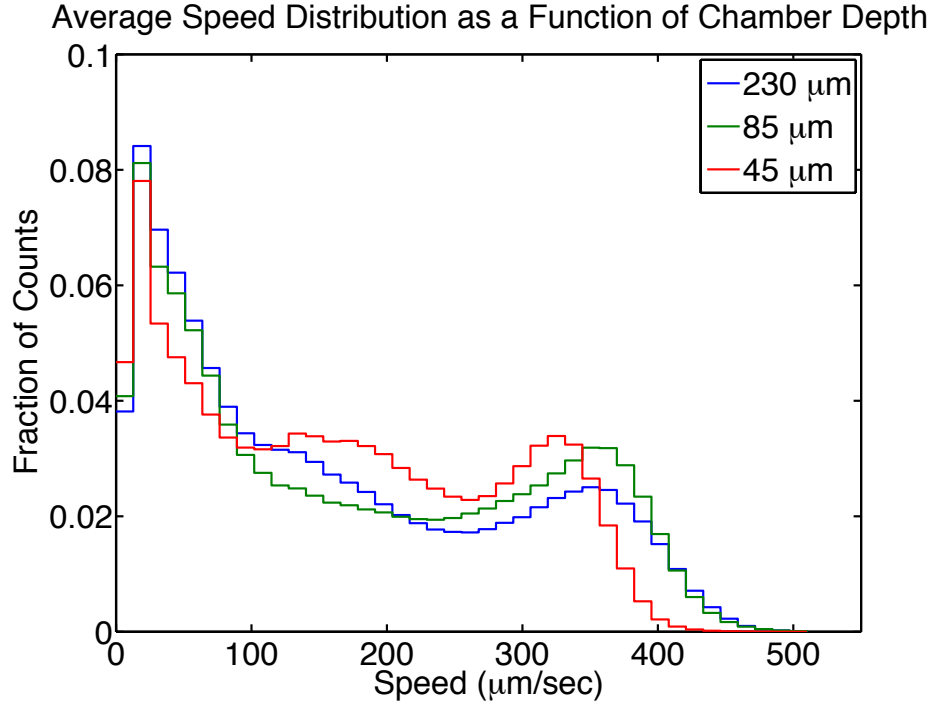


Figure A.1: **Averaged Histograms of Swimming Speed for Different Chamber Depths** - (N=15) for each Depth. The location of the high speed swimming peak changes when moving from 45 to 80  $\mu m$  , but not when the chamber depth is changed from 80  $\mu m$  to 230  $\mu m$  .

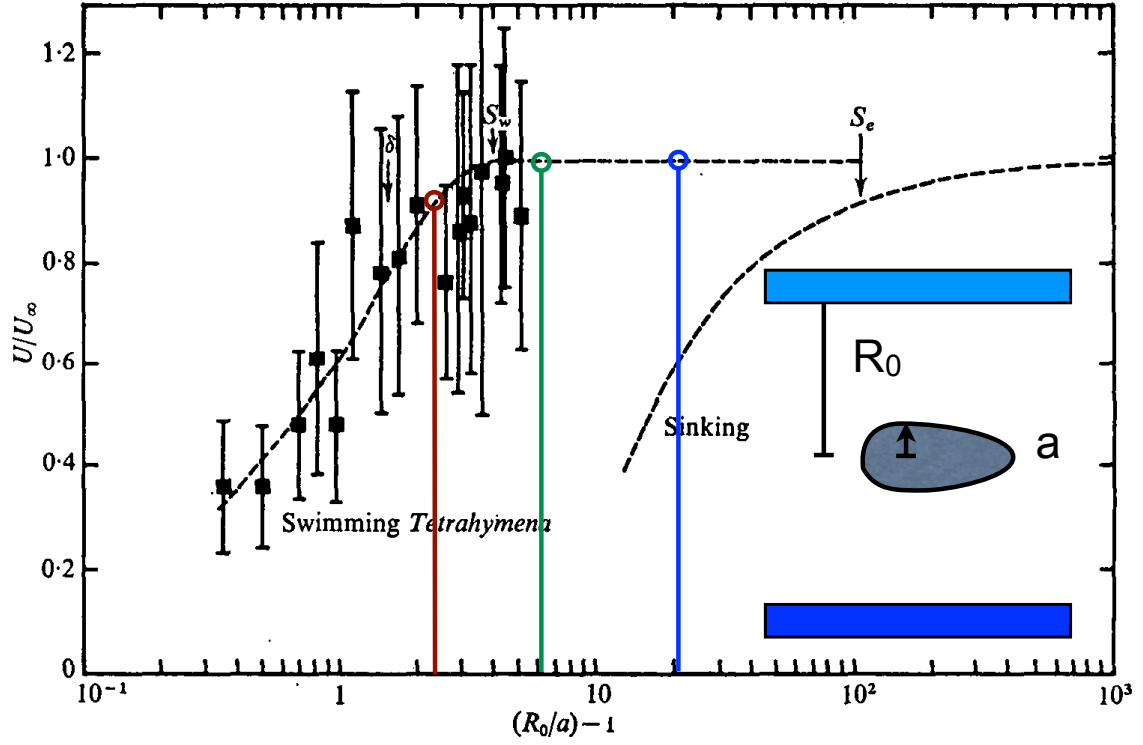


Figure A.2: **Ratio of Swimming Speeds in Proximity to a Wall as a fraction of Speed Far From the Wall** - The dashed lines represent the theoretical prediction and the Black Squares represent the experimental results of [50]. Our chosen chamber depths should put the wall drag in the regimes shown by the three stems, colors are as in Figure A.1, Red:  $45 \mu m$  , Green:  $80 \mu m$  , Blue:  $230 \mu m$  .  $R_0$  is the distance to the boundary and  $a$  is the ciliates radius.

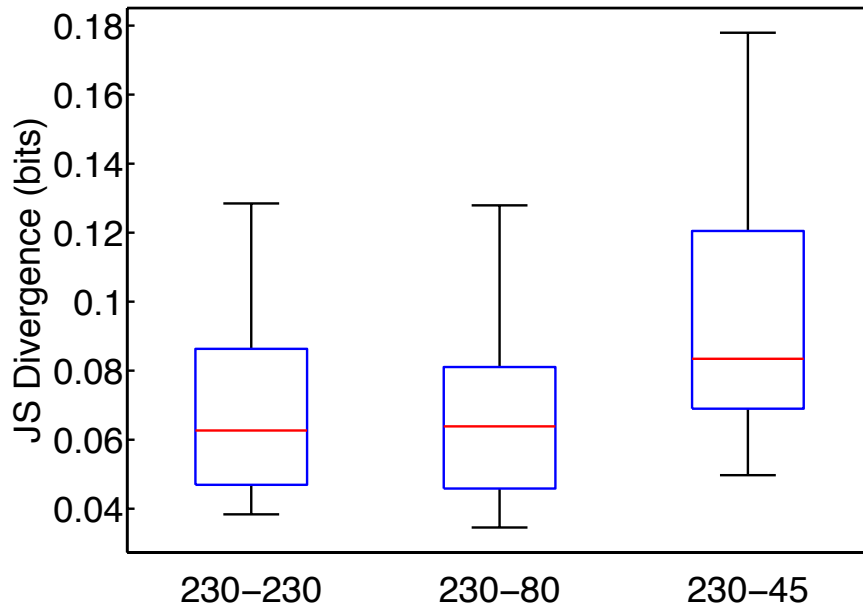


Figure A.3: **Distributions of JS divergence between Populations Assayed in Chambers of Different Depths** - (N=15) for each Depth. Red lines indicate medians, blu boxes show 25 and 75 percentiles, and whiskers indicate extrema. Swimming behavior is not significantly different for 230 and 80  $\mu m$  chambers, but shows measurable differences in 45  $\mu m$  chambers.

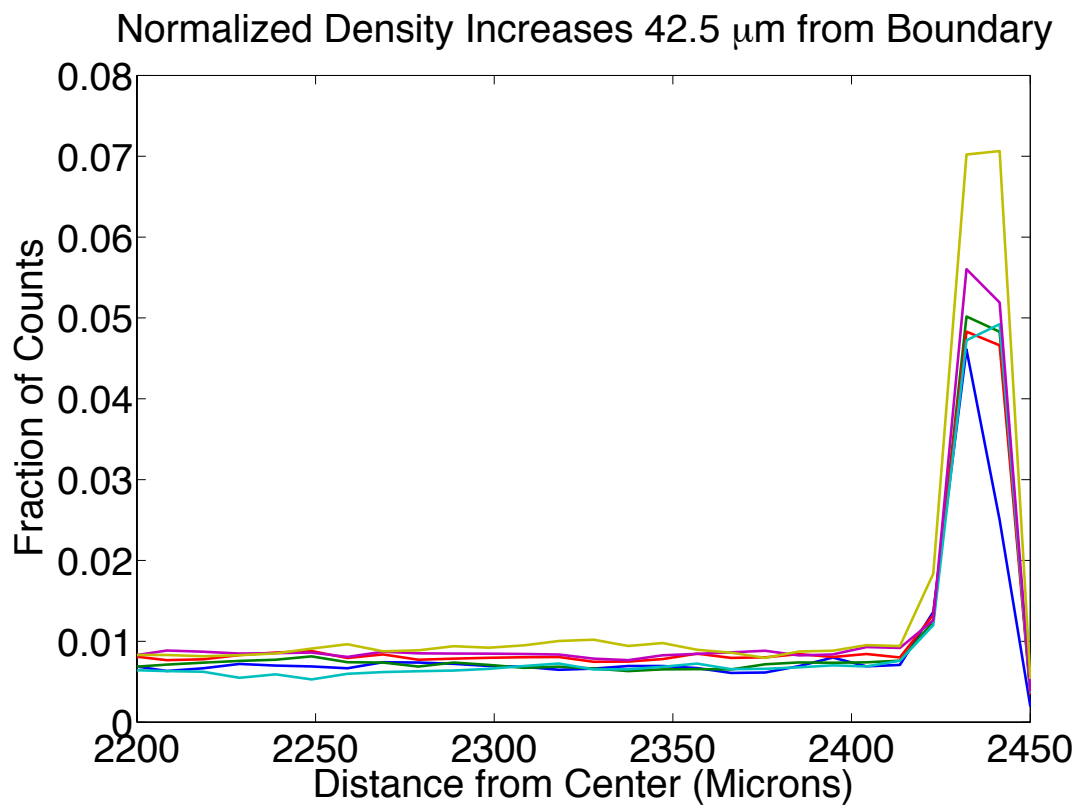


Figure A.4: **Radial distribution function near the wall** - for 6 individuals in the same chamber show the increase in normalized density near the wall.

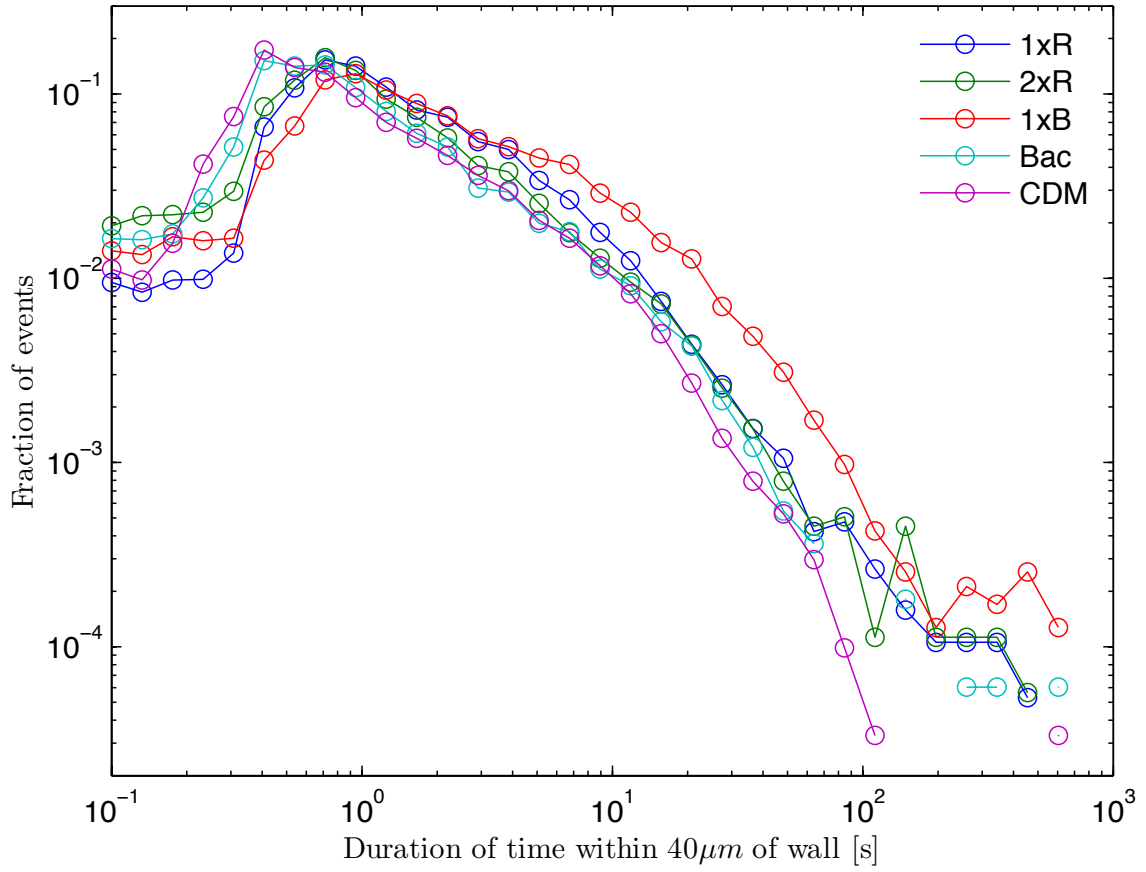


Figure A.5: **Distribution of Durations Spent Near the Wall** - for 5 Conditions, N=30 for each condition.

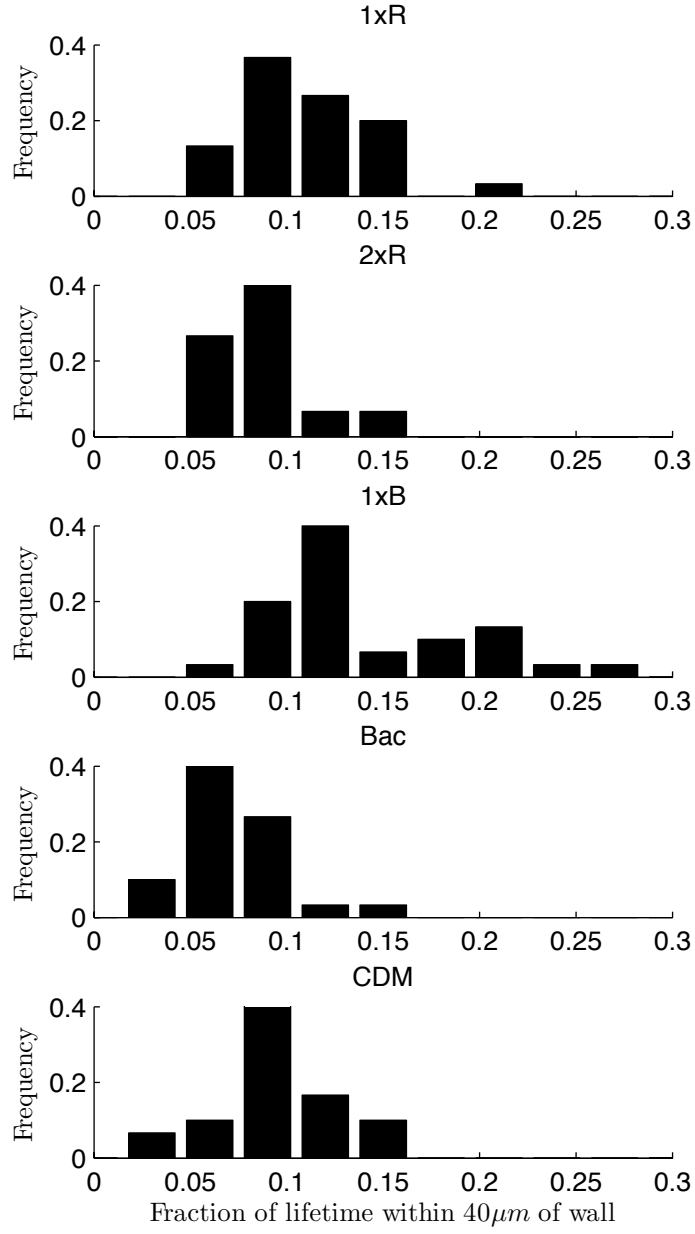


Figure A.6: **Fraction of Lifetime Spent Near the Wall** - for 5 Conditions,  $N=30$  for each condition.

With this definition of the wall/bulk boundary, we show that if an organism comes close to the wall and subsequently leaves, this interaction does not change the swimming behavior significantly compared to behaviors separated by the same amount of time in the bulk (Figure A.7).

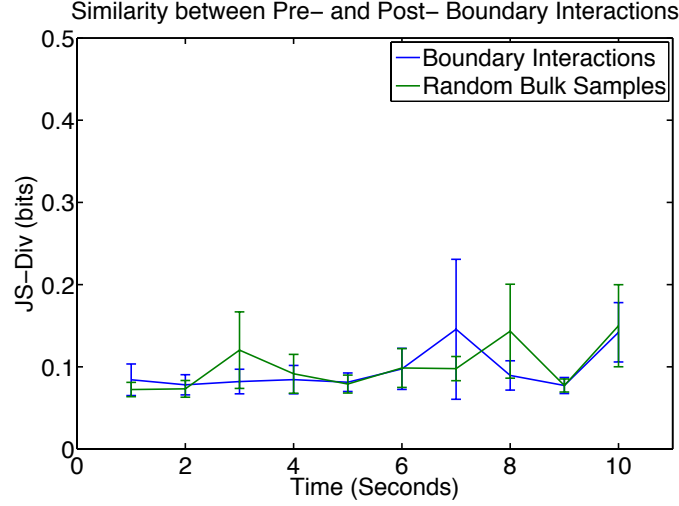


Figure A.7: **Behavioral Similarity for Behaviors Punctuated by an Interaction with the Wall** - Differences between behaviors separated by some amount of time  $t$  (x-axis) are similar, regardless of whether the organisms was near the wall (blue) or in the bulk (green) in the interceding interval.

### A.3 Chamber Isotropy

Chamber isotropy was investigated by looking at the spatial correlation between individuals in the same chamber, and comparing it to the spatial correlation of individuals in two different chambers. If there were long lasting spatial heterogeneities that resulting in some ‘preferred’ locations within the chamber, this would be evident as a greater spatial correlation between individuals in the same chamber. We observed

no difference between individuals in the same chamber when compared to individuals in different chambers, indicating no large-scale or long-lasting spatial heterogeneity.

The spatial correlation coefficient between individuals was calculated as

$$r = \frac{\sum_m \sum_n (\rho_{mn}^A - \bar{\rho}^A)(\rho_{mn}^B - \bar{\rho}^B)}{\sqrt{\sum_m \sum_n (\rho_{mn}^A - \bar{\rho}^A)^2 \sum_m \sum_n (\rho_{mn}^B - \bar{\rho}^B)^2}}$$

The distributions of spatial correlations between individuals in the same chamber are no different than that of those which lived in different chambers (Figure A.8). Some example individual density profiles are also shown, where darker represents higher density.

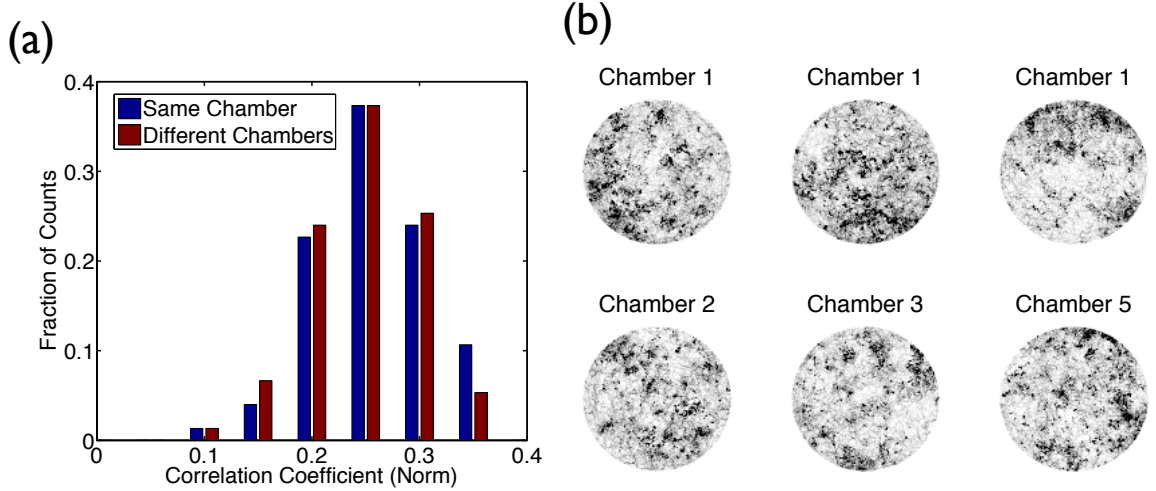


Figure A.8: **Spatial Correlation Between Individual Distribution Functions** - (a)  $N = 150$  (All Environments Considered). (b) shows some example density profiles from individuals in the same chamber (top row) and some individuals from different chambers (bottom row) in 1xR media.



## A.4 Nutrient Depletion

The volume a single *T. thermophila* is approximately one-millionth the volume of media in the chamber. This indicates that nutrient exhaustion should be negligible from over 3 generations. We have estimated the rate of nutrient uptake based on the work of [2] to be about  $10^{-12}g/min$ . Over the course of the experiment, we therefore estimate on the order of  $10^{-9}g$  to be consumed. The chamber starts with  $10^{-4}g$  of material, so this indicates that only 0.001% of the nutrients are used up over the course of the experiment.

# Appendix B

## System Details

### B.1 Temperature Control

The linear thermistor amplifier circuit was designed and built for this system by Seppe Kuehn (Figure B.1). The Peltier driver circuit was a custom designed Class B amplifier based on similar designs in [42].

### B.2 Chamber Fabrication

Chambers were constructed from PDMS using soft-lithography. SU8 - 2075 Negative photoresist spun onto a 4-inch diameter silicon wafer. 4 ml resist was dispensed onto the center of the wafer and the wafer was spun at 500 RPM for 10 seconds then at 1000 RPM for 30 seconds to a thickness of  $240\mu m$ . The soft bake was 7 minutes at  $65^{\circ}C$  followed by 45 minutes at  $95^{\circ}C$ . The wafer was then masked with a printed transparency mask and exposed to 575 mJ of radiation over 18 seconds and

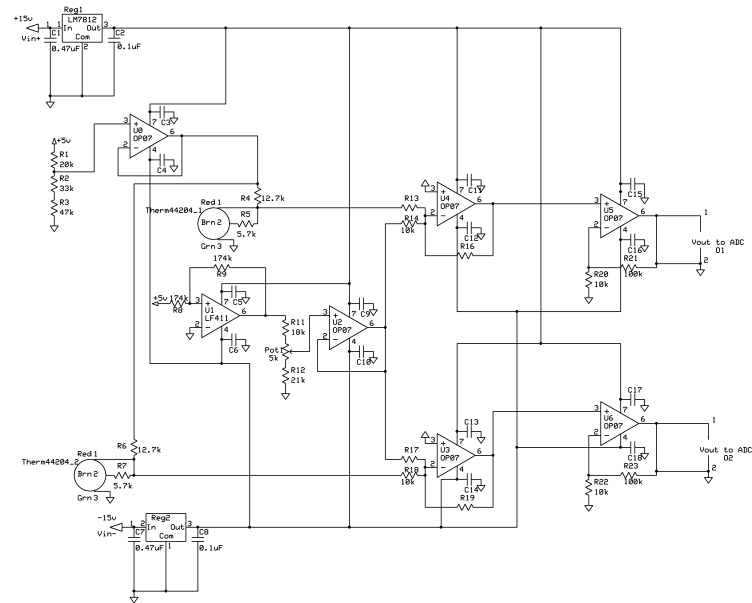


Figure B.1: **Circuit Diagram for Thermistor Amplifier** - Courtesy Seppe Kuehn.

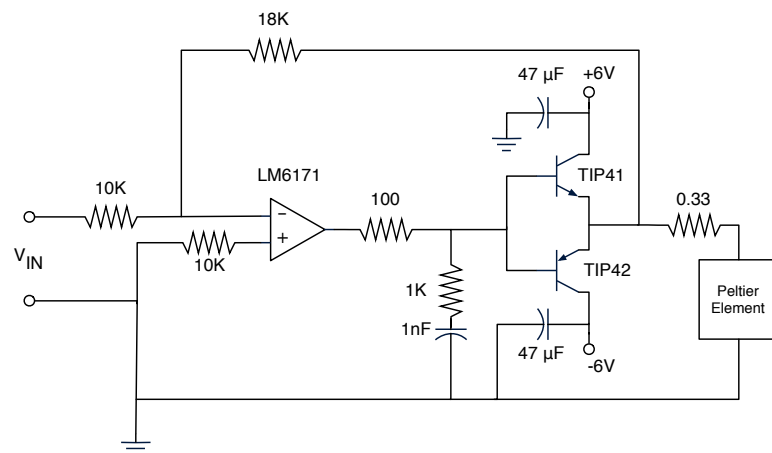


Figure B.2: **Circuit Diagram for Peltier Driver**

baked post exposure for 5 minute at 65°C and 15 minutes at 95°C. The resulting mold was developed for for 15 minutes in SU-8 developer. This mold was exposed to tridecafluoro - 1,1,2,2 tetrahydrooctyl trichlorosilane for 1 hour mild vacuum to prevent sticking of PDMS. Sylguard 184 Silicone Elastomer (Dow Corning, Midland MI), was mixed with Sylguard 184 curing agent at a ratio of 8:1 and poured onto the silicone mold. This was cured for 25 minutes at 80°C then allowed to cool. The PDMS chamber was trimmed using a razor blade and holes were punched using a syringe needle. The chamber was then plasma treated and sealed to a glass coverslip before being cured for an additional hour at 80°C.

### B.3 Detection

Each image is captured and imported into Matlab (The MathWorks, Natick MA) using the VideoIO toolbox (Gerald Daley) as a matrix  $I(t) = \{p_{ij}(t) | 1 \leq i \leq 1200, 1 \leq j \leq 1600\}$  where  $p_{ij}(t)$  represents the value of the intensity of the pixel at location  $(i, j)$  at time  $t$ . A time delayed projection is calculated as a pseudo-background image for detecting moving objects. For a time delay of  $\tau$ , the background image is given as

$$B_{ij} = \max(p_{ij}(t_0), p_{ij}(t_0 + \tau))$$

A motion detection matrix is then constructed as

$$M(t)_{t=t_0:\tau} = I(t) - B$$

Connected pixels in this matrix for which the pixel intensity value  $m_{ij}(t)$  is greater than some threshold  $T$  represent any objects that are not in the same location at times  $t_0$  and  $\tau$ . These objects are filtered to retain only the expected number of objects  $N_0$  (with some exceptions) in the expected size range of the *Tetrahymena* ( $> 25 \text{ pixels}^2$ ). This is done by ordering the objects  $x_n$  from largest to smallest and saving first  $k$  objects where

$$k = \max(N_0, |\{x : \text{area}(x) > A_{\max}\}|)$$

Where  $|\bullet|$  indicates the cardinality of a set and  $A_{\max}$  is the cutoff size for artifacts, in this case  $60 \text{ pixels}^2$ . For each object  $x_n$  for the  $n = 1 : k$  retained objects, the centroid, area, orientation and eccentricity are recorded using the *regionprops* function in the Matlab Image Processing Toolbox (The MathWorks, Natick MA).

## B.4 Cost Functions

Trajectories are created by matching objects in frame  $t$  with those in  $t + 1$ . This problem can be posed as a linear assignment problem. Each object  $x_i$  in frame  $t$  can be linked to an object  $x_j$  in frame  $t + 1$  for a cost  $c_{\text{link}}(x_i, x_j)$ . The cost function used is the sum of the euclidian distance between objects and a weighted difference in area. Alternatively each  $x_i$  can remain unlinked for a cost  $c_{\text{lose}}$  and each  $x_j$  for a cost  $c_{\text{find}}$ . We use the hungarian algorithm to find the assignment matrix  $\hat{A}_{ij}$  which has the minimal total cost  $C$  where  $C = \sum_i \sum_j \hat{A}_{ij} c_{ij}$ . In our formulation, the cost

to lose or find a particle was set to 1.25 times the maximum of linked assignment costs in all previous frames, this was determined empirically. Links are assigned using the hungarian algorithm. Objects are first linked pairwise frame by frame into segments using this cost function. These segments are aggregated and linked in a second linear assignment matching using a cost function which incorporates euclidian distance and change in median segment area, but also includes the time gap and the vector change in velocity between the end of one segment and the beginning of another. The trajectories formed from the joined segments are then checked for fidelity (see Appendix D ).

The list of  $N$  detected particles in frame  $t$  is given as  $X^N(t) = \{x_n(t)\}$  for  $n = 1 : N$ . Trajectories would ideally be created by finding the minimum cost for one-to-one assignments over the bipartite graph with  $X^N(t)$  and  $X^M(t+1)$  as vertices and  $N = M$ . We seek to find  $\hat{A}_{ij}$  which minimizes the sum of costs

$$\sum_i \sum_j A_{ij} C_{ij}$$

This is know as the linear assignment problem. Because particles may disappear and artifacts may arise during object detection,  $N$  does not always equal  $M$ . We have used the cost matrix formulation of [25] to account for these possibilities. Thus the cost matrix  $C$  is constructed from the actual  $N \times M$  linking matrix  $c_{ij} = c_{link}(x_i(t), x_j(t+1))$  in the following way.

$$C = \begin{vmatrix} c_{ij} & c_{lose}(X^N(t)) * I_N \\ c_{find}(X^M(t+1)) * I_M & c_{ij}^T \end{vmatrix}$$

where  $c_{lose}$  is the cost of an object in frame  $t$  remaining unassigned,  $c_{find}$  is the cost of an object in frame  $t+1$  remaining unassigned, and  $I_n$  is the  $n \times n$  identity matrix. The  $c_{ij}^T$  term is added to satisfy the requirements of the linear assignment formulation. We have used

$$cost_{lose} = cost_{find} = 1.25 * \max(\{\hat{A}_{ij}(t_0) * C(t_0) | t_0 = 1 : t, i < N, j < M\})$$

.

where  $N = |X^N(t)|$  and  $M = |X^M(t+1)|$ , and  $i < N$  and  $j < M$  ensures only linking assignments are considered.

The frame to frame linking cost function is given by

$$c_{link}(x_i(t), x_j(t+1)) = \Delta r(x_i(t), x_j(t+1)) + \alpha \Delta A(x_i(t), x_j(t+1))$$

where  $\Delta r$  is the displacement of the centroids of the objects and  $\Delta A$  is the difference in area.  $\alpha$  is a scaling factor. In our case,  $\alpha$  is set to 0.2. For the segment linking step, let  $S^N$  be the set of  $N$  segments  $s_i$  to be joined. The cost function to join segment  $s_i$  to  $s_j$ ,  $i \neq j$  is,

$$c_{link}(s_i, s_j) = \Delta r_{ij}(s_i, s_j) + \tau \Delta t_{ij}(s_i, s_j) - \Delta v_{ij}$$

where  $\Delta r_{ij}$  is the displacement from the end of segment  $i$  to the beginning of segment  $j$ ,  $\Delta t_{ij}$  is the time of the gap in frames,  $\tau$  is a scaling factor (0..5 in this case), and  $\Delta v_{ij}$  is the projection of the velocity in the last frame of  $s_i$  onto the velocity in the first frame of  $s_j$ , normalized to  $\Delta r_{ij}$ .



# Appendix C

## Behavioral Analysis

A behavior is a sequence of observable actions. Humans learn a variety of actions and string them together into behaviors. For example, series of actions steps, hops, and leaps - result in behaviors such as walking, running, or skipping. The quantitative study of behavior requires a method to measure behavioral differences. While running is clearly different from skipping, it is not a priori clear how to measure such a difference. Classifying behaviors such as walking, running or skipping is called scoring and the classes are called stereotyped behaviors. Scoring, however, can become difficult when boundaries between different classes of behaviors become unclear and when environmental or genetic changes alter existing classes or cause new ones to emerge. Thus, inspecting the action sequences directly, rather than scoring can be in general more appropriate. In what follows we will introduce statistical methods that describe behaviors using only their underlying actions and thus avoiding scoring.

## C.1 Comparing behaviors through distributions divergence

Following our example we will denote the behaviors *walk*, *run*, and *skip* as sequences of actions  $(step)_n$ ,  $(leap)_n$ , and  $(leap, hop)_n$  respectively. Intuitively, run is more similar to skip than to walk, because they share a  $(leap)$  action. We will extend this intuition to construct a measure of distance between behaviors based on how frequently each of the observed actions is performed. Two behaviors with few actions in common will have a large distance relative to two that have many common actions.

Histograms quantify frequencies of observed actions and provide a standardized output across behaviors. Since histograms are estimates of probability distributions, a measure that compares distributions, such as the Kullback-Leibler divergence, is appropriate. We have chosen the related Jensen-Shannon divergence, which satisfies the requirements of a metric, one of which is symmetry. Symmetry is important because the similarity between two behaviors should not depend on the order in which they are listed. If  $p_i$  and  $q_i$  are the frequencies of action  $i$  for two behaviors, then the JS divergence between those behaviors is given by the following equation:

$$D(p, q) = \frac{1}{2} \left[ \sum_i p_i \log \left( \frac{p_i}{m_i} \right) + \sum_i q_i \log \left( \frac{q_i}{m_i} \right) \right]$$

With this metric, the distance between two behaviors that have no actions in common is one and the distance between two behaviors that have the same actions, identical in frequency, is zero. Behaviors with more actions in common are less distant.

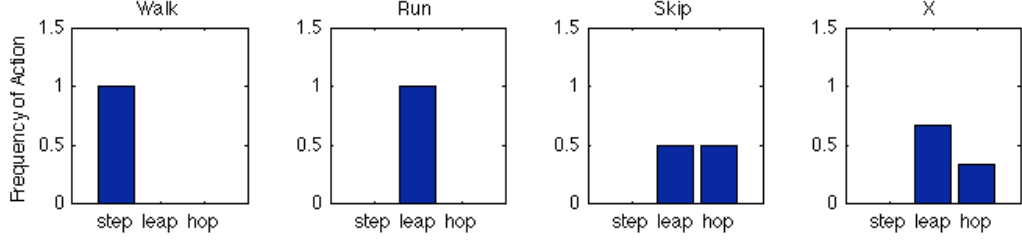


Figure C.1: **Histograms of Actions for a Simple Example** These histograms show the frequency of each action for our simple examples of behavior. They can be represented as vectors, for example,  $walk=[1\ 0\ 0]$  and  $skip=[0\ 0.5\ 0.5]$ . Because walk does not share any actions with the other behaviors, it has a distance of 1 from everything. skip and X are the most similar with a distance of 0.02. run is closer to X ( $D = 0.19$ ) than it is to skip ( $D = 0.31$ )

Figure C.1 demonstrates the use of this metric to quantify behavioral distance in our simple example. The distance between run and skip by this metric is 0.3113, while the distance between run and walk is 1. Figure C.1(d) extends our example. Let us imagine a novel behavior is observed,  $(leap, leap, hop)_n$ . We can immediately note the weakness of scoring, as there is no stereotyped behavior to describe this. However, using our formalism, we can immediately compare it quantitatively to the existing behaviors.

The example shown in Figure C.1 is illustrative of how one can compare behaviors using their underlying action sequences and highlights how stereotyping can fail when novel behaviors arise. However, in this example, the actions themselves are stereotyped. This can obscure important differences that might result from differing external stimuli or internal states. The motions that comprise a step might differ if

that step is taken by an old passerby or if it is taken by an Olympic speed walker. In what follows, we will define actions using quantitative descriptions of the motions themselves, avoiding all stereotyping.

Locomotion behaviors can be quantified in many ways, from the acceleration of limbs measured in horse gait analysis, to the large-scale movements of birds tracked using GPS tags. These measurements result in a time series of a kinematic variable. This variable could for example be the angular speed of a knee joint or the velocity of a bird in flight. To illustrate our statistical method of quantifying behavior, let us take a hypothetical measurement of some quantity (X), relevant to behavior. In our example above, this variable could be linear speed, with *hop* being the slowest, *step* being intermediate, and *leap* being the fastest. Just as we constructed histograms to represent distributions of the discrete stereotyped actions, we can construct histograms of the variable X. We can use these histograms to quantify behavioral differences.

## **C.2 Changeability: variations of behavior during lifetime**

The ability to measure behavioral differences allows us to examine how an individual's behavior changes during in time. The way a person moves changes with age and varies through different stages of life. If we record a sequence of actions, we can divide these sequences into intervals and compute the frequency of the actions observed in each interval. This defines a behavior for each time interval. The behaviors at different

intervals can then be compared using the Jensen-Shannon metric as described above.

This procedure is shown for our hypothetical variable  $X$  in Figure C.2.

The individual in Figure C.2(a) exhibits behavior that changes in time, moving between lower and higher levels of the relevant variable (denoted by  $X$ ) for the first two thirds of the measurement before fluctuating wildly for the remainder. While stereotyping may appear appropriate for early times, where behavior seems to fall into one of two stereotyped behaviors, it is unreasonable to classify the later behavior into either of these stereotypes. Figure C.2 illustrates how we can measure the behavioral changes of an individual during in time. We call the matrix describing these differences an individuals changeability matrix.

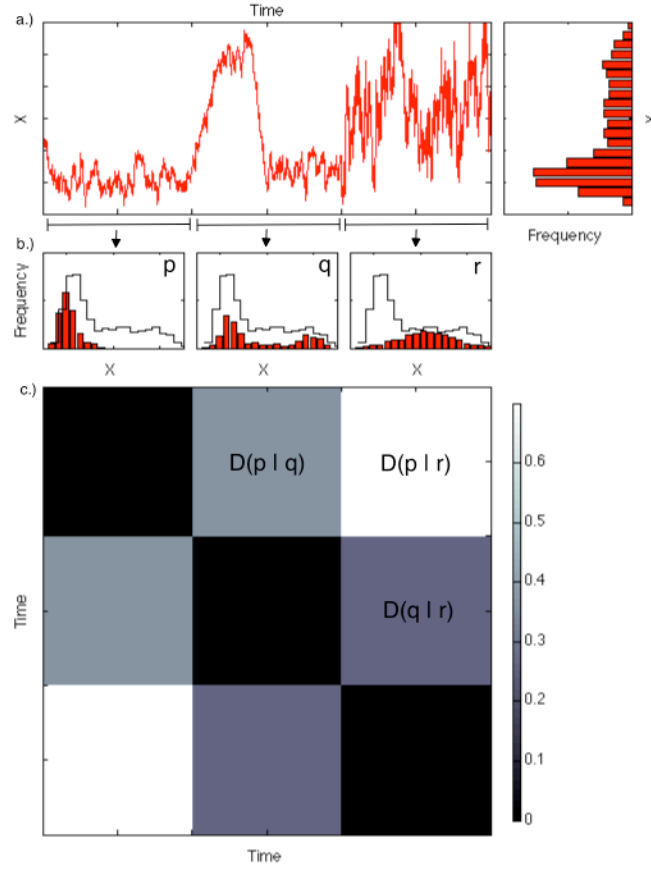


Figure C.2: **Measuring Individual Behavioral Changes** Consider a measurement, in time, of a parameter relevant to behavior ( $X$ ), (see text). (a) A time series of this parameter for an individual. The right panel shows a histogram of the entire time series. (b) Histograms for three non-overlapping intervals of the time series shown in (a). These intervals are labeled  $p$ ,  $q$ , and  $r$  sequentially. The black line gives the histogram for the entire time series as shown in (a), right panel. (c) A symmetric matrix of similarity measurements for each interval (b) with each other interval in (b). We term this matrix the changeability matrix since it quantifies the changes in an individual's behavior in time.

## C.3 Individuality: Differences of behavior among individuals

In addition to computing changeability, we can use this formalism to compare the behavior one individual to that of others, a measure which we call individuality. When comparing individuals, it can be simpler to compare their average of their behaviors, rather than the each of behaviors separately. While it is unclear how to average stereotyped behaviors, our approach allows averaged behaviors to be represented as the average of their respective histograms. Thus, the single histogram generated from the sequence of all observed actions represents an individuals average behavior and the JS divergence between two such histograms is the distance between those individuals. With such histograms, we can compare the average behaviors between individuals and look for differences resulting from genetic or environmental changes. This is shown for three individuals and our hypothetical variable X, in Figure C.3.

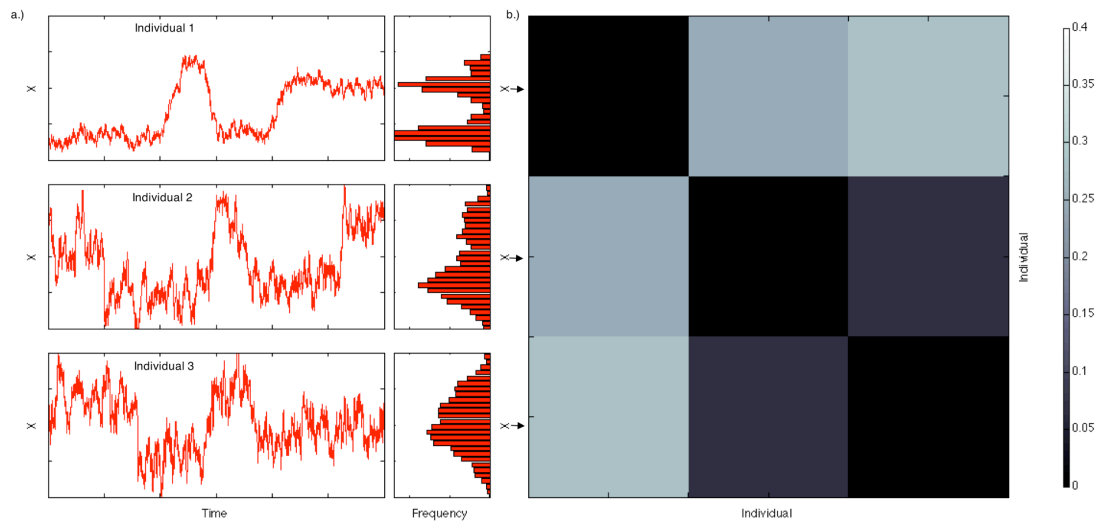


Figure C.3: **Measuring Differences Between Individuals** The Individuality matrix measures differences in behavior between individuals. (a) A time series of parameter ( $X$ ) for three individuals. Histograms for these three time series are shown at the right. (b) An individuality matrix for the three individuals shown in (a)



# Appendix D

## Tracking, Error and Uncertainty

### D.1 Tracking Fidelity

Correct assignments are essential to maintain the identity of individuals over the course of the experiment. In general, when the ratio of the frame-to-frame displacement and the inter object distance is 1 and each is small ( $< 212\mu m$ ) there is the possibility of an erroneous assignment. Sets of contiguous frames in which this ratio is close to one will be called "crossover events". This ratio is a function of the average speed of the objects, the frame rate of the video, and the density of objects. In our apparatus, we can easily maintain identity with a density up to 8 individuals. In an average experiment, there are 25 events per object pair per "lifetime" in which the above ratio is close to one (less than 1.6). Across all experiments, 96.2% of the crossover events are correctly assigned by the automated tracking system. All crossover events are inspected manually and those that are not correctly assigned are corrected by hand. The validity of the hand scored trajectory is verified by detailed

analyses of the movie and comparison of a variety of parameters including the resulting lifetimes of individuals, the speed of each, the angular acceleration, and the size of each before and after the crossing.

## D.2 Imaging Uncertainty

Uncertainty in the position as a result of the imaging hardware (e.g. illumination variation, pixel noise, optical point-spread function) was evaluated by recording stationary objects of known size using the USAF test target (Edmund Scientific). Video was recorded for 1 minute and the resulting 1000 frames were processed using the same custom MATLAB algorithms developed for tracking (Mathworks, Natick MA). The centroid position of the objects from the resulting segmented images were recorded and their deviations across the 1000 frames were measured to be on average 0.2 microns, with a maximum of 1 micron. There was no clear dependence on this variance with the size of the objects between 10 and 4500 square pixels.

## D.3 Divergence Estimators

Given histograms  $\mathbf{H} = \{\mathbf{h}^m\}$  which are estimates of distributions  $\mathbf{P} = \{\mathbf{p}^m\}$ , we would like to quantify the bias and the uncertainty of our estimates of the divergence  $D(\mathbf{p}^1|\mathbf{p}^2)$ . Detailed derivations of the analytic expressions for the bias and uncertainty are given in [17].

### D.3.1 Bias

The bias is given as the difference between the expected value of  $D(\mathbf{h}^1|\mathbf{h}^2)$  and  $D(\mathbf{p}^1|\mathbf{p}^2)$ . Following [17], let us consider two sequences of observations of length  $N$ , each of which can be in  $i = 1..B$  ‘states’ according to the unobserved distributions  $\mathbf{p}^1$  and  $\mathbf{p}^2$ . In our work, these states are bins in our histograms. Let  $n_i$  represent the number of observations in each state  $i$ . Thus our histograms  $\mathbf{h}_i$  are given by  $\frac{n_i}{N}$ .

$$bias = \langle D(\mathbf{h}^1|\mathbf{h}^2) \rangle - D(\mathbf{p}^1|\mathbf{p}^2)$$

The bias is independent of the underlying distributions  $\mathbf{P}$ , and depends only on the ratio of the number of occupied states  $B^*$  to the number of observations  $N$ . State  $i$  is occupied if  $\mathbf{p}_i \neq 0$ . The bias is then given by,

$$\langle D(\mathbf{h}^1|\mathbf{h}^2) \rangle - D(\mathbf{p}^1|\mathbf{p}^2) = \frac{B^* - 1}{2N \ln(2)}$$

The bias represents the average divergence you would expect to measure between independent samples drawn from an identical unobserved distribution. This value is systematically biased upwards from its true value due to sampling error, and decreases with better sampling (larger  $N$ ). In this work, the bias was corrected analytically for each pair of histograms using the appropriate  $B^*$  values. Generally, the bias is about 1% for divergences associated with sub-lifetime changeability measurements ( $D(p_t^n|p_{t'}^n)$ ), and negligible (1000 times smaller) for full lifetime individuality divergences ( $D(p^n|p^m)$ ).

### D.3.2 Uncertainty

The uncertainty is given as the variance of  $D(\mathbf{h}^1|\mathbf{h}^2)$ .

$$uncertainty = V[D(\mathbf{h}^1|\mathbf{h}^2)]$$

[17] shows that the uncertainty  $V[D(\mathbf{h}^1|\mathbf{h}^2)]$  depends on the number of observations  $N$ , and depends only on terms of  $O(1/N^2)$  and smaller. We have performed bootstrap estimates of the uncertainty and find that it does indeed scale as  $O(1/N^2)$  and is on the order of  $10^{-6}$  bits, negligible for both the changeability and for the individuality data presented.

# Appendix E

## Changeability and Individuality Distributions

These figures summarize the changeability and individuality distributions for each of the five conditions in which wild type *T. thermophila* was assayed. Each individual's changeability can be represented by a matrix which has  $(T/t)^2$  elements, where  $T$  is the length of the trajectory and  $t$  is the length of the window. This matrix is symmetric, so we can collect the non-redundant entries from the upper triangle. For each individual, a histogram of the  $[(T/t)^2 - (T/t)]/2$  unique entries is constructed and the median is calculated and distribution of these medians for a population of 30 individuals is plotted as a box and whisker plot (Figure E.1). For individuality matrices, the distributions of the upper triangle of the individuality matrix  $I(N, M) = (P^n|P^M)$  is shown as a box and whisker plot (Figure E.2)

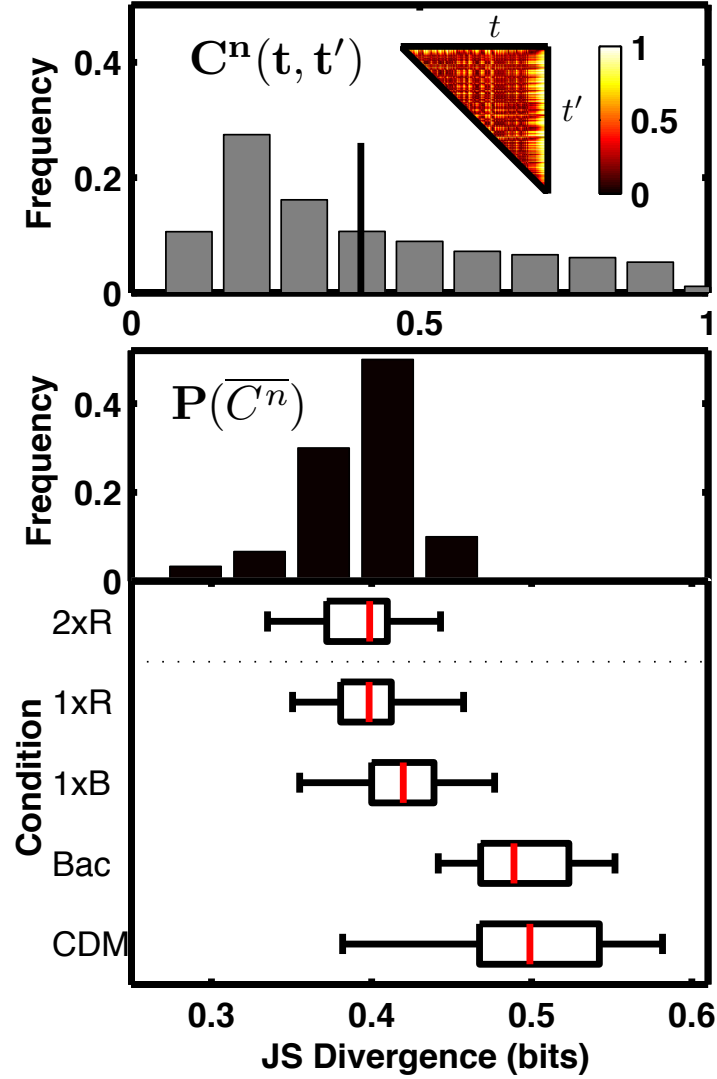


Figure E.1: **Changeability Distributions for All Environments for *T. thermophila*** a) The  $C^N(t, t')$  matrix is symmetric so only the upper triangle is shown (inset), along with a histogram of all non-redundant entries of  $C^N(t, t')$  (grey bars). This histogram reflects the magnitude of the differences in behaviors exhibited by an individual in its lifetime. The mean of this distribution, denoted  $\overline{C^N(t, t')}$  and shown by the vertical black line in (a) we call the mean changeability for individual  $n$ . (b upper panel) Shows the distribution of  $\overline{C^N(t, t')}$  for all 30 individuals in a single experimental condition, 2x SPP (2xR) in this case. We denote this by  $P(\overline{C^N})$ . (lower panel) Box and whisker plots of for all five experimental conditions studied here. The boundaries of the box and the red line show the quartiles of the distribution and the median respectively. The whiskers give the 5th and 95th percentile of the data.

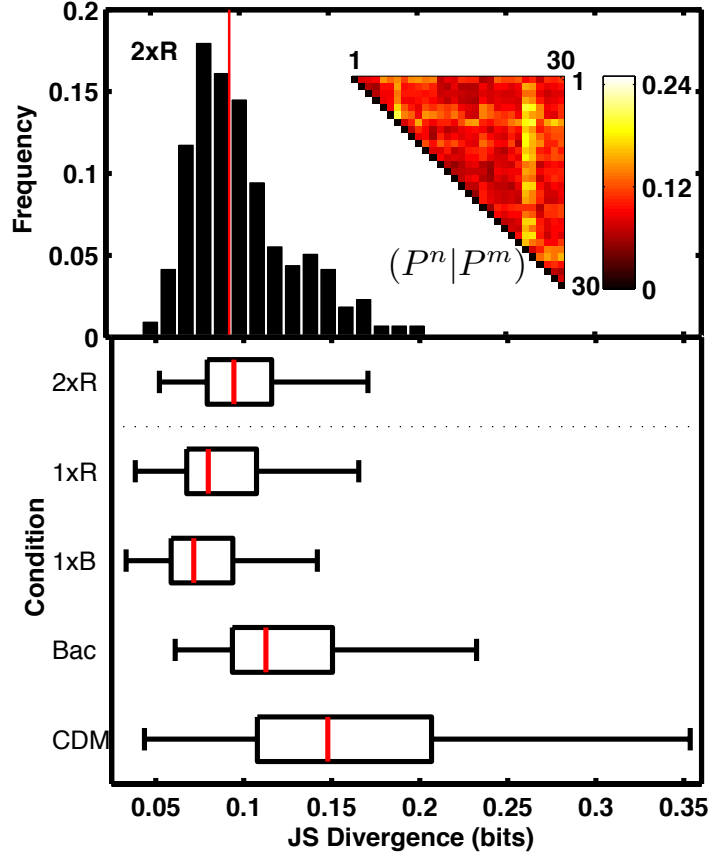


Figure E.2: **Individuality Distributions for All Environments for *T. thermophila*** For all of the data presented here we compare histograms from the entire lifetime of each individual. (a) Individuality is measured by  $I(N, M) = (P^N | P^M)$ , or the pair-wise difference in behavior between all individuals in a given condition.  $I(N, M)$  is a symmetric 30 x 30 matrix for each condition. The upper triangle of  $I(N, M)$  is shown for 2xR (inset). A histogram of all non-redundant entries of  $I(N, M)$  reflects the distribution of behavioral differences between individuals (black bars). The vertical red line indicates the median. (b) Box and whisker plot as in figure 2(b) except for individuality in each experimental condition.

# Appendix F

## Protocols

### F.1 *Tetrahymena* sp. Information

Table F.1: *Tetrahymena* sp. Strain Information - Strains were obtained from the Cornell Tetrahymena Stock Center <http://tetrahymena.vet.cornell.edu>

Abbrev.	ID	Species	Strain	Collected By	Location (Lat Lon)
Tb	SD01609	<i>T. borealis</i>	—	—	—
Tt:VT	SD01566	<i>T. thermophila</i>	20488-4	Paul Doerder	42 50.802 -72 42.596
Tt:PA	SD01554	<i>T. thermophila</i>	19869-1	Paul Doerder	41 38.414 -79 54.691
Tt:NH	SD01554	<i>T. thermophila</i>	20469-4	Paul Doerder	42 57.763 -72 07.451

### F.2 Cell Culture

Cells were maintained long term in soybean culture and passaged every 6 months in duplicate. Soybean stocks were prepared by autoclaving 10  $cm^3$  of distilled water with a soybean. This was inoculated with 50  $mm^3$  of an exponentially growing



*Tetrahymena sp.* culture [45]. Exponentially growing cultures for both passaging soybean stocks and performing experiments were prepared in the same way. 50  $cm^3$  of the growth media (1xR for passaging cultures) was prepared and autoclaved, and inoculated with 50  $mm^3$  from a soybean stock. This culture was allowed to grow for approximately 48 hours at 23C without shaking.

# Bibliography

- [1] Julius Adler. Chemotaxis in bacteria. *Science*, 153:708–716, August 1966.
- [2] Anders Andersen and Per Hellung-Larsen. Division competence in Tetrahymena: Determination of minimum cell volume and rate of nutrient uptake. *Journal of Cellular Biochemistry*, 41, 1989.
- [3] Natalie Balaban, Jack Merrin, Remy Chait, Lukasz Kowalik, and Stanislas Leibler. Bacterial Persistence as a Phenotypic Switch. *Science*, 305:1622–1625, September 2004.
- [4] Andres Bendesky, Makoto Tsunozaki, Matthew V Rockman, Leonid Kruglyak, and Cornelia I Bargmann. Catecholamine receptor polymorphisms affect decision-making in *C. elegans*. *Nature*, 472:313–318, April 2011.
- [5] H.C. Berg and L. Turner. Cells of *Escherichia coli* swim either end forward. *Proceedings of the National Academy of Sciences USA*, 92:477, 1995.
- [6] Howard Berg and Douglas Brown. Chemotaxis in *Escherichia coli* analysed by three-dimensional tracking. *Nature*, 239:500–504, 1972.
- [7] Waheb Bishara, Serhan O Isikman, and Aydogan Ozcan. Lensfree optofluidic microscopy and tomography. *Annals of Biomedical Engineering*, 40:251–262, February 2012.
- [8] E.H. Blackburn and C.W. Greider. Identification of a specific telomere terminal transferase activity in Tetrahymena extracts. *Cell*, 43:405–413, 1985.
- [9] Kristin Branson, Alice A Robie, John Bender, P Perona, and MH Dickinson. High-throughput ethomics in large groups of *Drosophila*. *Nature Methods*, 6:451–457, 2009.
- [10] J.E. Brownell, J. Zhou, T. Ranalli, R. Kobayashi, D.G. Edmondson, S.Y. Roth, and C.D. Allis. Tetrahymena histone acetyltransferase A: a homolog to yeast Gcn5p linking histone acetylation to gene activation. *Cell*, 84:843–851, 1996.
- [11] ME Cates, D. Marenduzzo, I. Pagonabarraga, and J. Tailleur. Arrested phase separation in reproducing bacteria creates a generic route to pattern formation. *Proceedings of the National Academy of Sciences USA*, 107:11715, 2010.

- [12] F Paul Doerder and Clifford Brunk. *Natural Populations and Inbred Strains of Tetrahymena*, volume 109. Elsevier Inc., March 2012.
- [13] Marie-Anne Félix and Christian Braendle. The natural history of *Caenorhabditis elegans*. *Current Biology*, 20:R965–9, 2010.
- [14] Tom Fenchel. Microbial behavior in a heterogeneous world. *Science*, 296:1068, 2002.
- [15] J Frankel, L M Jenkins, and L E DeBault. Causal relations among cell cycle processes in *Tetrahymena pyriformis*. An analysis employing temperature-sensitive mutants. *Journal of Cell Biology*, 71:242–260, October 1976.
- [16] Wei Geng, Pamela Cosman, Charles C Berry, Zhaoyang Feng, and William R Schafer. Automatic tracking, feature extraction and classification of *C. elegans* phenotypes. *IEEE Transactions on Bio-medical Engineering*, 51:1811–1820, 2004.
- [17] I Grosse, P Bernaola-Galván, P Carpena, R Román-Roldán, J Oliver, and HE Stanley. Analysis of symbolic sequences using the Jensen-Shannon divergence. *Physical Review E*, 65:41905, 2002.
- [18] S Gueron, K Levit-Gurevich, N Liron, and J J Blum. Cilia internal mechanism and metachronal coordination as the result of hydrodynamical coupling. *Proceedings of the National Academy of Sciences USA*, 94:6001–6006, June 1997.
- [19] A.P. Gupta and RC Lewontin. A study of reaction norms in natural populations of *Drosophila pseudoobscura*. *Evolution*, 36:934–948, 1982.
- [20] E.M. Hedgecock and R.L. Russell. Normal and mutant thermotaxis in the nematode *Caenorhabditis elegans*. *Proceedings of the National Academy of Sciences USA*, 72:4061, 1975.
- [21] John J Hopfield. Physics, computation, and why biology looks so different. *Journal of Theoretical Biology*, 171:53–60, 1994.
- [22] S Elizabeth Hulme, Sergey S Shevkoplyas, Alison P McGuigan, Javier Apfeld, Walter Fontana, and George M Whitesides. Lifespan-on-a-chip: microfluidic chambers for performing lifelong observation of *C. elegans*. *Lab On A Chip*, 10:589–597, 2010.
- [23] Konstantin G Iliadi and Gabrielle L Boulianne. Age-related behavioral changes in *Drosophila*. *Annals of the New York Academy of Sciences*, 1197:9–18, June 2010.
- [24] JD Jacobs and JC Wingfield. Endocrine control of life-cycle stages: A constraint on response to the environment? In *Condor*, pages 35–51. Univ Washington, Dept Zool, Seattle, WA 98195 USA, 2000.

- [25] Khuloud Jaqaman, Dinah Loerke, Marcel Mettlen, Hirotaka Kuwata, Sergio Grinstein, Sandra L Schmid, and Gaudenz Danuser. Robust single-particle tracking in live-cell time-lapse sequences. *Nature Methods*, 5:695–702, 2008.
- [26] Ronald Konopka and Seymour Benzer. Clock mutants of *Drosophila melanogaster*. *Proceedings of the National Academy of Sciences USA*, 68:2112–2116, September 1971.
- [27] Ekaterina Korobkova, Thierry Emonet, Jose Vilar, Thomas Shimizu, and Philippe Cluzel. From molecular noise to behavioural variability in a single bacterium. *Nature*, 428:574–578, 2004.
- [28] K. Kruger, P.J. Grabowski, A.J. Zaug, J. Sands, D.E. Gottschling, and T.R. Cech. Self-splicing RNA: autoexcision and autocyclization of the ribosomal RNA intervening sequence of *Tetrahymena*. *Cell*, 31:147–157, 1982.
- [29] I Richard Lapidus. Growth and Division Kinetics of Asymmetrically Dividing *Tetrahymena thermophila*. *Journal of Theoretical Biology*, 106:135–140, 1984.
- [30] V Leick and P Hellung-Larsen. Chemosensory behaviour of *Tetrahymena*. *BioEssays*, 14:61–66, January 1992.
- [31] J. Lin. Divergence measures based on the Shannon entropy. *IEEE Transactions on Information Theory*, 37:145–151, 1991.
- [32] Taejin L Min, Patrick J Mears, Lon M Chubiz, Christopher V Rao, Ido Golding, and Yann R Chemla. High-resolution, long-term characterization of bacterial motility using optical tweezers. *Nature Methods*, 6:831–835, October 2009.
- [33] D L Nanney. Corticotype transmission in *Tetrahymena*. *Genetics*, 54:955–968, October 1966.
- [34] E M Nelsen. Transformation in *Tetrahymena thermophila*. Development of an inducible phenotype. *Developmental Biology*, 66:17–31, September 1978.
- [35] E M Nelsen, J Frankel, and L M Jenkins. Non-genic inheritance of cellular handedness. *Development*, 105:447–456, March 1989.
- [36] T. Piersma and J. Drent. Phenotypic flexibility and the evolution of organismal design. *Trends in Ecology and Evolution*, 18:228–233, 2003.
- [37] Stephen Quake and Axel Scherer. From micro-to nanofabrication with soft materials. *Science*, 2000.
- [38] L Rasmussen and L Modeweg-Hansen. Cell multiplication in *Tetrahymena* cultures after addition of particulate material. *Journal of Cell Science*, 12:275, 1973.

- [39] Hanna Salman and Albert Libchaber. A concentration-dependent switch in the bacterial response to temperature. *Nature Cell Biology*, 9:1098–1100, September 2007.
- [40] Robert Schleif, Winand Hess, Solomon Finkelstein, and D. Ellis. Induction kinetics of the L-arabinose operon of *Escherichia coli*. *Journal of Bacteriology*, 115:9–14, 1973.
- [41] Matthew Scott and Terence Hwa. Bacterial growth laws and their applications. *Current Opinion in Biotechnology*, 22:559–565, August 2011.
- [42] Paul Sherz. *Practical Electronics for Inventors 2/E*. McGraw-Hill/Tab Electronics, November 2006.
- [43] John L Spudich and Daniel E Koshland. Non-genetic individuality: chance in the single cell. *Nature*, 1976.
- [44] R Stavis and Hirschberg R. Phototaxis in *Chlamydomonas reinhardtii*. *The Journal of cell biology*, 59:367–377, 1973.
- [45] Melody Sweet and C. David Allis. Long-Term Storage of Unfrozen *Tetrahymena* Cultures in Soybean Medium. *Cold Spring Harbor Protocols*, 2006:pdb.prot4472–pdb.prot4472, August 2006.
- [46] L. Szablewski, P Andreassen, A. Tiedtke, Florin Christensen, Florin Christensen, and L Rasmussen. *Tetrahymena thermophila*: Growth in synthetic nutrient medium in the presence and absence of glucose. *Journal of Eukaryotic Microbiology*, 38:62–65, 1991.
- [47] J W Szostak and E.H. Blackburn. Cloning yeast telomeres on linear plasmid vectors. *Cell*, 29:245–255, May 1982.
- [48] J B Tenenbaum, V de Silva, and J C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290:2319–2323, December 2000.
- [49] Jose Vilar, Calin Guet, and Stanislas Leibler. Modeling network dynamics: the lac operon, a case study. *The Journal of Cell Biology*, 161:471–476, May 2003.
- [50] Howard Winet. Wall drag on free-moving ciliated micro-organisms. *Journal of Experimental Biology*, 59:753–766, 1973.
- [51] Fred W Wolf, Aylin R Rodan, Linus T-Y Tsai, and Ulrike Heberlein. High-resolution analysis of ethanol-induced locomotor stimulation in *Drosophila*. *The Journal of Neuroscience*, 22:11035–11044, 2002.