

1 Inter-individual Variation in Circadian Rhythm Across Species of Halictid Bees:
2 Characterizing Varying Clocks and Motif Classification
3

4 By

5
6 Sofía Meléndez Cartagena

7
8 A thesis submitted to the Department of Biology
9

10
11 FACULTY OF NATURAL SCIENCES
12 UNIVERSITY OF PUERTO RICO
13 RÍO PIEDRAS CAMPUS
14

15
16
17 In partial fulfillment of the requirements for the degree of
18 Master of Science
19

20
21 May 26 2021
22

23 Rio Piedras, Puerto Rico
24

26 Thesis Approval

27

28

29

30 FACULTY OF NATURAL SCIENCES

31 UNIVERSITY OF PUERTO RICO

32 RÍO PIEDRAS CAMPUS

33

34 In partial fulfillment of the requirements for the degree of

35

36 Master of Science

37

38 Thesis Committee:

39

40 _____ Advisor

41 José L. Agosto Rivera, PhD

42

43 _____ Co-Mentor

44 Patricia Ordoñez, PhD

45

46 _____

47 James Ackerman, PhD

48

49 _____

50 Steven E Massey, PhD

Dedication

To all the hands that made this work move forward. This is not only a labor of love, but also of community. To the curious citizens that saw me out in the field collecting and offered their help and wisdom. To my colleagues at the graduate program that kept me sane. To my friends and family, for whom my ramblings about my work must sound like gibberish, but encourage me to keep talking anyway. And of course, to my science enthusiast parents, for helping me grow into a better scientist.

Acknowledgments

I would like to thank my mentors Dr. José L. Agosto Rivera and Dr. Patricia Ordoñez for taking me in and giving me the opportunity to grow as a scientist and as a person. To Dr. James Ackerman, who has seen me grow from a baby scientist into a not-quite-grown scientist. To Dr. Steven E Massey, for asking me the difficult questions and leading me to think about science from a different perspective. I would also like to thank Elvia Melendez-Ackerman for adopting me as an honorary lab member and offering her time, wisdom, and resources.

To my peers from the graduate program, especially Jonathan A Lopez Colon, Diego A. Rosado Tristani and Elif Kardas, without your help, field work would not have been the same.

A special big thanks to my friends, Melanie E. Martines Arrocho, Nadira Yusif Rodriguez, Brian X. Leon Ricardo, and Alicia Figueroa Carlo. You did everything from reading my drafts, helping me practice my talks and even going out in the field with me. Your support is invaluable and I'm lucky to have you as friends.

Lastly, I would like to thank my undergraduate mentees, Alejandro Armas, Luis Roman Lizasoain, Beatriz A. Morales Grimany, Milexis A Santos Vega, and Andrea V. Velez Velez. I loved growing with you and wish you happiness and satisfaction in your future endeavors.

97	Table of Contents	
98		
99	Thesis Approval	1
100	Dedication	2
101	Acknowledgments	3
102	Table of Contents	4
103	List of Tables and Figures	5
104	Abbreviations List	7
105	Abstract	9
106	Author Biography	10
107	Introduction	11
108	Chapter 1	19
109	Chapter 2	47
110	Discussion/Conclusion	92
111	Literature Cited	97
112	Appendix	107
113		
114		

115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137

List of Tables and Figures

Introduction

Table 1: Levels of sociality as adapted from Michener, C D (1969).

Chapter 1

Figure 1: Habitat (A-B) and Species Observed (C-F).

Figure 2: Female *S.curvicornis* exhibit short period phenotype under constant darkness (<24 h endogenous circadian rhythm)

Figure 3: A summary of the variations in the circadian rhythm as observed in:

A) *Lasioglossum malachurum*, B) *Lasioglossum enatum* and *Lasioglossum ferreirii*

Figure 4: *L.malachurum* exhibits a variety of circadian phenotypes under constant dark conditions.

Figure 5: Description of the circadian behaviors exhibited by *L. ferreirii* under constant dark conditions.

Figure 6: Description of the circadian behaviors exhibited by *L. enatum* under constant dark conditions.

Figure 7: Summary of descriptive and inferential statistics

Chapter 2

Figure 1: Locomotor activity over time data from *L. malachurum* displays heterogeneity.

Figure 2: Average of all 98 individuals in the dataset is not representative of any one individual due to heterogeneity.

Figure 3: A mock-up of a normalized time series transformed with SAX.

Figure 4: A mock-up of how KNN works.

Figure 5: Mock-up illustrating Decision Trees (left) and Random Forest (right).

Figure 6: PAM with $K = 4$ for the Lomb Scargle Periodogram resulted in highly representative discrete clusters.

Figure 7: Clusters resulting from L. S are separated by the shape of the periodogram Part 1.

Figure 8: Clusters resulting from L. S are separated by the shape of the periodogram Part 2.

Figure 9: PAM with $K = 3$ for Autocorrelation coefficient clusters of low representation power.

Figure 10: PAM with $K = 3$ for Autocorrelation coefficient clusters by thickness of Autocorrelation.

Figure 11: K-means with $K = 3$ for Average Daily Activity clusters by frequency.

Table 1: Results of KNN by varying PAA Size (PAA) and Alphabet Size (Alphabet). Sliding Window (SW) was continually equal to 48.

Table 2: Results Trees by varying PAA Size (PAA) and Alphabet size (Alphabet). Sliding Window (SW) was continually equal to 48.

Figure 12: Prevalence in NAs is due to poor consistent classification.

Appendix

Figure A.1: Methods for capturing and husbandry of the bees

161

162

Abbreviations List

163

Abbreviation	Meaning
1. <i>L.malachurum</i>	<i>Lasioglossum malachurum</i>
2. <i>S. curvicornis</i>	<i>Systropha curvicornis</i>
3. <i>L. ferrerii</i>	<i>Lasioglossum ferrerii</i>
4. <i>L. enatum</i>	<i>Lasioglossum enatum</i>
5. <i>L.parvus</i>	<i>Lasioglossum parvus</i>
6. <i>L. gemmatum</i>	<i>Lasioglossum gemmatum</i>
7. LAM	Locomotion Activity Monitor
8. RS	Rhythm Strength
9. R	Rhythmic
10.WR	Weakly Rhythmic
11.A	Arrhythmic
12.L.S	Lomb Scargle
13.M.L	Machine Learning

14. SAX	Symbolic Aggregate approXimation
15. KNN	K Nearest Neighbors
16. PAM	Partitioning Around Medoids
17. BoW	Bag of Words
18. TF-IDF	Term frequency–inverse document frequency
19. PCA	Principal Component Analysis
20. SW	Size of Word
21. PAA	Piecewise Aggregate Approximation
22. DT	Decision Trees
23. RF	Random Forest
24. TP	True Positive
25. FP	False Positive
26. FN	False Negatives

164

165

166

167

Abstract

How sociality interacts with other behaviours is a long standing question in insect biology. Simultaneously, in chronobiology, there is an unanswered question concerning how sociality influences patterns of daily activity. Past studies have established the viability of utilizing hymenopterans to describe variable circadian behaviour. Here, we intend to take a step further and establish Halictid bees as a model for cross-species comparisons of circadian rhythms. To this effect, we describe four species of Halictid bees and compare the variability of their internal clocks. We found that variability in circadian rhythms parallel complexity in sociality. We also created a machine learning pipeline to facilitate describing heterogeneous locomotor activity data. The computational experiments showed that the bee locomotion dataset transformed by Symbolic Aggregate approXimation and classified by decision trees yielded the best results. With our findings, we expect to set the basis for finding the true influence of sociality on biological clocks.

Author Biography

Sofía Meléndez Cartagena is interested in learning about the evolutionary pressures that select for sociality. In her current work, she is evaluating how sociality may affect circadian rhythms, and is developing methods to best access them. Her pursuits span anywhere from doing extensive field work to spending long hours in front of the computer building Machine Learning models. When not doing science, Sofía splits her time between enjoying the theater of the mind and doing activism work.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240

Introduction

Introduction

Circadian rhythms are believed to be a universal mechanism of life (Helm and Visser 2010). They determine when flowers open, when animals are hungry, and even when the temperature of a mammal's body might be at their highest or their lowest. It has been stated that these biological clocks follow a circadian rhythm, which is approximately synchronized to Earth's 24-hour rotation using signals from the environment (Roenneberg, T., et.al 2003).

Drosophila melanogaster is the main model organism used for the study of circadian rhythms. Thanks to its controlled genetic stocks created for laboratories, we have been able to describe biological clocks at the molecular level (Dubowy C & Sehgal A. 2017). Nonetheless, the genetic homogeneity and ease of manipulation that makes fruit flies attractive models leaves them lacking when it comes to answering questions about individual differences in behavior. Learning how biological clocks work in organisms that deviate from the canon may allow us to answer questions that are currently left unanswered by the *Drosophila* model. One such question is on how sociality may affect daily activity patterns.

The most studied model by far for cross-comparative studies in reference to sociality are Hymenopterans (ants, wasps and bees) and Isopterans (Termites) (West-Eberhard 2003). Although they are not usually kept in captivity and their genetics are not as well-understood as *Drosophila*'s, the diversity of social behaviours is well-documented. Furthermore, past phylogenetic studies have found that the genes responsible for the circadian circuitry in the brain as found in Hymenopterans are more similar to those of mammals in contrast to *Drosophila*

(Rubin, et. al 2006 and Ingram, et. al 2012). Their well-known behaviors and similarity with mammals makes Hymenopterans a model of interest for transferability and comparison for mammalian behaviours.

Sociality in insects is defined as a spectrum, and, depending on the study, one might find that there are a myriad of definitions. The two biggest factors that influence those definitions are based on how an individual of a species may relate to other members of the same species, and the adaptations in behaviors concerning reproduction and brood care (Michener 1969, Toth and Rehan 2017). Here, we focus on four distinct levels of sociality: solitary, communal, primitively eusocial, and facultatively eusocial. These levels of sociality are individually defined in Table 1. To clarify, by facultatively eusocial insects we refer to primitively or strongly eusocial species whose life history suggest that they were originally solitary, or have evolved a solitary life cycle, or possess an adaptable developmental system that can express social or solitary behavior (Eickwort 1996).

Table 1: Levels of sociality as adapted from Michener (1969). 1 means present and 0 means absent.

	Caste and Division of Labor	Overlapping Generations	Cooperative Work on Cells	Structurally similar reproductive female caste (if present) may survive alone	Progressive feeding
Solitary	0	NA	0	1	0
Communal	0	0	0	1	1
Primitively eusocial	1	1	1	1	0/1
Strongly eusocial	1	1	1	0	0/1

Previous studies such as (Moore et al. 1998, Bloch et al. 2001, Giannoni-Guzman et al. 2020) have already begun establishing the use of honey bees as non-canonical subjects of study in circadian science. For example, Moore et al. 1998 and Bloch et al. 2001 established that with age, the way circadian rhythms are expressed in bees change. Furthermore, in Giannoni-Guzman et al. 2020, it was found that even within the same age group and apparent caste, there is evidence of shift work, which influences the daily activity patterns of individuals within a hive and reflects in interindividual variation. The honey bee as a subject of study proves to be a powerful tool to answer questions that are too difficult to explore under the constraints of the typical *Drosophila* model. However, if one wishes to study the influence of sociality in circadian rhythms, *Apis mellifera* is a lacking model as due to

297 apparent lack of diversity, genus-level *Apis* are exclusively eusocial. One would have
298 to compare them to members of different genera to extract any sort of meaningful
299 conclusions on how varying levels of sociality may affect the expression of circadian
300 phenotypes, as was done in Giannoni-Guzman et al. 2014. Ultimately, the optimal
301 way to observe the relationship between sociality and circadian rhythms would be to
302 observe a set of species with high levels of social plasticity.

303 One such group of bees exist, the tribe Halictini, which in the past has been
304 suggested as a subject of study for the interaction of socio-comparative studies
305 (Schwarz et al. 2007, Bloch and Grozinger 2011, Toth and Rehan 2017). In the first
306 chapter of this work, we describe and compare four different species of halictid bees
307 whose social organization encompasses the spectrum of sociality. *Systropha*
308 *curvicornis* (Scopoli) is a solitary species (Patiny and Michez 2007; Patiny et al.
309 2008; Danforth et al. 2008) and *Lasioglossum malachurum* (Kirby) (*L. malachurum*)
310 is an obligately eusocial species (Richards 2000; Wyman and Richards 2003) from
311 Greece. *Lasioglossum (Dialictus) ferrerii* (Baker) is communal. *Lasioglossum*
312 (*Dialictus*) *enatum* (Gibbs) has not been actively studied, but it is related to
313 primitively eusocial species (Eickwort 1988), both of these last bees are from Puerto
314 Rico. The bees from Greece were captured the same day and same time while
315 visiting the flowers of *Campanula arvensis*, and in addition, were kept in the same
316 environmental conditions in the laboratory. Similarly, the bees from Puerto Rico were
317 captured at the same day and same time while they visited *Momordica charantia*,
318 *Sida acuta*, and *Bidens alba*.

319 We found that although *L. enatum* and *L. ferreiri* highly related and share the
320 same environment, there exist distinct differences in their behavior. Furthermore, as
321 an overall observation, bees that are not solitary at the intraspecies level express
322 individual differences in circadian behavior. Which begs the question: what causes
323 these differences, if not the environment nor their species? Identifying the root
324 causes of these intraspecies behavioral differences may give us insight into how
325 biological clocks are entrained by non-environmental cues.

326 The individual differences exhibited by *L. malachurum* galvanized an
327 exploration towards a preprocessing in the traditional circadian analysis pipeline.
328 Because these bees have such diverse expressions of circadian phenotypes, we
329 found ourselves dividing the bees into discrete groups and describing the population
330 by the use of subpopulations. This process was the manual equivalent of grouping
331 individuals using machine learning. Once we were confident that the groups
332 observed within the population of *L. malachurum* were discrete and informative for
333 circadian purposes, we took the next step and designed a machine learning pipeline.

334 In chapter 2, we explore the use of said Machine Learning (M.L) pipeline.
335 Machine learning is a subset of Artificial Intelligence. It is defined as the process of
336 solving real-life problems by gathering data sets and using algorithms to create
337 models based on those data sets (Burkov, 2019). Although machine learning has
338 many applications, we employed it for two distinct purposes: clustering and
339 classification. Clustering is a type of unsupervised machine learning where
340 algorithms are not given predetermined categories (labels) to group data into. The
341 data, based on measurements of similarity, instead groups together naturally (Géron,

2018). Classification is a type of supervised learning where data already has the desired solutions (labels). A classification algorithm learns patterns from a labeled dataset and uses this information to classify new data into these same categories (Géron, 2018).

The use of M.L in circadian science is usually relegated to the evaluation of genes and proteins (Agostinelli et. al 2016 and Anafi et. al 2014). It is our understanding that the use of M.L as a preprocessing step to locomotor analysis is a novel application of this tool, which presents the challenge of having to design a pipeline and workflow from scratch.

Using common transformations for circadian analyses (Lomb Scargle Periodogram, Average of Daily Locomotor Activity and Autocorrelation) (Refinetti et. al 2007) we explored the natural grouping of the various transformations of the *L. malachurum* dataset. The approach was univariate, where we applied PAM or K-means to just one of the transformations. This culminated in three different clustering results, one per transformation. None of the three transformations resulted in groups that were of circadian significance, and prompted us to explore the use of classification. We divided the data for *L. malachurum* into three broad categories, and then transformed the time series data using Symbolic Aggregate approXimation (SAX) for ease of analysis and dimensionality reduction. The data was transformed with different combinations of SAX parameters and were classified using K Nearest Neighbors, Decision Trees and Random Forest classification algorithms. Lastly, we compared the performance of all three algorithms. Decision Trees achieved the best results, followed closely by K Nearest Neighbors.

We hypothesize that sociality is one of the keys to understanding the complexities of circadian rhythms. In this work, we meet our set goal by; First, describing four species of halictid bees with varying degrees of sociality. Second, by developing a preprocessing step using machine learning that is easily incorporated into the traditional circadian pipeline. This preprocessing step facilitates the grouping of individuals of *L. malachurum* into discrete groups once trained. While it stands to be seen if the process is transferable to other species, its effectiveness can nevertheless prove useful in future studies where *L. malachurum* is involved. All of the caveats notwithstanding, our findings can serve as the basis for a larger body of work that elucidates the true relationship between sociality and circadian rhythms.

390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412

Chapter 1

Relationship Between Inter-individual Variation in Circadian Rhythm and Sociality: A
case Study Using Halictid Bees

**Work done in collaboration with: Carlos A. Ortiz-Alvarado, Claudia S. Cordero-
Martínez, Alexandria F. Ambrose, Luis A Roman Lizasoain, Milexis A Santos
Vega, Andrea V Velez Velez, Jason Gibbs, Theodora Petanidou, Thomas
Tscheulin, John T. Barthell, Victor H. González, Tugrul Giray, José L. Agosto-
Rivera**

Relationship Between Inter-individual Variation in Circadian Rhythm and Sociality: A case Study Using Halictid Bees

Abstract:

The bee family Halictidae is considered to be an optimal model for the study of social evolution due to its remarkable range of social behaviors. Past studies in circadian rhythms suggest that social species may express more diversity in circadian behaviors than solitary species. However, these older studies did not make appropriate taxonomic comparisons. To further explore the link between circadian rhythms and sociality, we examine four halictid species with different degrees of sociality, three social species of *Lasioglossum*, one from Greece and two from Puerto Rico, and a solitary species of *Systropha* from Greece. Based on our previous observations, we hypothesized that species with greater degree of sociality will show greater inter-individual variation in circadian rhythms than solitary species. We observed distinct differences in their circadian behavior that parallel differences across sociality, where the most social species expressed the highest inter-individual variation. We predict that circadian rhythm differences will be informative of sociality across organisms.

Introduction:

Understanding the evolutionary link between solitary and eusocial lineages, and their adaptive behaviors, such as those expressed in reproduction and brood care, is a perennial question in insect evolutionary biology (Toth and Rehan 2017). Insects, and in particular hymenopterans, have been useful in observing how

sociality is related to other types of behaviors, such as competitor effects (Peters et al. 2017). A potential behavior to be evaluated is circadian rhythms, as it has been proposed to be governed by demands arising from sociality not only in insects but also in mammals (Mistlberger 2004; Giannoni-Guzman et al. 2014; Beer and Helfrich-Förster 2020).

Circadian rhythms can be viewed as a biological clock that most living organisms possess. These biological clocks regulate processes such as gene expression, behavior, body temperature, and sleep-wake patterns. Biological clocks follow a rhythm that is approximately synchronized to Earth's 24-hour rotation using signals from the environment called zeitgebers. This process of synchronization needs active reestablishment, and it is called entrainment (Roenneberg et al. 2003).

Circadian rhythms have been studied in a wide range of organisms, from plants, invertebrates, birds, and mammals (Helm and Visser 2010). The traditional model animal to study this phenomenon is the fruit fly, *Drosophila melanogaster*, where biological clocks are described at the molecular level (Dubowy and Sehgal 2017). Although this model has been pivotal to the understanding of circadian rhythms, the lack of genetic diversity in the fruit fly reduces the relevance of the model because it limits questions regarding individual differences in rhythms.

Past studies, such as those of Bloch et al. (2001) and Moore et al. (1998), revealed that the rhythmicity of honey bees changed with age. Additionally, Giannoni-Guzman et al. (2020) showed that foragers in the wild display discrete categories that suggest temporal shift work. An earlier study from Giannoni-Guzman

and colleagues (2014) compared the endogenous period of three different variants of honey bees (*Apis mellifera carnica*, *Apis mellifera caucasica* and *Apis mellifera* gAHB) as well as similarly sized insects from different orders and families. They found that honey bees and paper wasps (*Polistes crinitus* and *Mischocyttarus phthisicus*) had a larger degree of circadian period variation within the population in comparison to *D. melanogaster*. They mentioned several possible explanations for their observations, one of them being sociality. In a more recent study, Beer and Helfrich-Förster (2020) explore this connection further and note that the development of the circadian circuitry varies between an eusocial (*Apis mellifera*) and a solitary species (*Osmia bicornis*). In particular, they observe that eusocial individuals are born with an undeveloped circadian clock while the solitary individuals emerge with it fully developed, and attribute these differences to their opposite levels of sociality. However, because these two past studies were done with species spanning from different taxonomic groups, it would be difficult to support their claims without taking phylogeny into account. Nevertheless, these works do give a basis to ask how circadian rhythms vary and are an integral part of the survival strategy and organization of these animals. Moreover, it leads us to consider that the level of sociality in different organisms may play a role in their daily activity patterns. Specifically, one could suppose that complexity in levels of sociality of an insect may be reflected in individual differences in circadian rhythms of individuals of the same population.

Halictidae (Hymenoptera) is a bee family considered to be a great model for the study of social evolution due to its exceptional diversity in respect to social

behavior within and among species and populations (Schwarz et al. 2007). *Lasioglossum Curtis* is one of the two genera in the tribe *Halictini* that displays eusocial behavior, but also includes solitary representatives and a range of intermediate social categories (Danforth et al. 2003; Gibbs et al. 2012). Additionally, past studies have shown plasticity in the social behavior even among populations of the same species (Eickwort et al. 1996; Field 1996; Field et al. 2010; Richards et al. 2003; Soucy and Danforth 2002; Richards 2000). Depending on environmental conditions such as elevation, latitude and seasonality, halictid bees might display different modes of sociality. Species with social nests may revert to solitary behavior at high latitudes and altitudes (Eickwort et al. 1996; Packer et al. 1983; Field et al. 2010) or based on access to mates (Yanega 1988/1989). Jeanson et al. (2008) studied members of a solitary species, *Lasioglossum (Ctenonomia)* sp. NDA-1 and observed the results of having them nest in pairs. They observed that after some time together, the individuals in the nest started to show signs of division of labor. This plasticity and diversity of behavior, in addition to the close taxonomic relation, makes *Halictidae* an optimal model for observing the relationship between sociality and circadian rhythm (Bloch and Grozinger 2011).

To better understand how social behaviors can be associated with circadian rhythms in insects, we have set out to document the rhythm in four species of halictid bees that span a gradient of social complexity. *Systropha curvicornis* (Scopoli) (*Halictidae: Rophitinae*), a solitary pollinator specialist (Grozđanić and Mučalica 1966) considered ancestrally solitary within the family Halictidae (Patiny and Michez 2007; Patiny et al. 2008; Danforth et al. 2008), and three species of *Lasioglossum* (L.

ferreri, *L. enatum* and *L. malachurum*), which were selected because of their varying levels of social behavior (Eickwort 1988; Wyman and Richards 2003; Gibbs 2018). Species of *Lasioglossum* likely had a common ancestor capable of eusocial nesting but have reverted multiple times to other levels of sociality (Danforth et al. 2003; Brady et al. 2006; Gibbs et al. 2012). *Lasioglossum* (*Dialictus*) *ferrerii* (Baker) and *L.* (*Dialictus*) *enatum* (Gibbs) is found in the Caribbean, whereas *L. malachurum* is found across Europe; the first species nests communally, that is, each individual contributes to nest construction and reproduction (Michener 1974; Eickwort 1988). Although *L. enatum* has not been thoroughly studied, this species is part of a species complex that includes weakly eusocial species. Namely, *L. gemmatum* (Smith) and *L. parvum* (Cresson) from Jamaica and the Bahamas, which exhibit reproductive division of labour (Eickwort 1988). There is no evidence of morphologically defined castes beyond reproductive status in these two species, and thus we assumed this is likely the case for *L. enatum* in Puerto Rico, where we conducted our experiments. In contrast, *L. malachurum* (Kirby) is an obligately eusocial species with a morphologically well-defined queen and worker castes (Richards 2000; Wyman and Richards 2003). *Lasioglossum malachurum* is known to display varying degrees of behaviors depending on location (Richards 2000). In Lesbos, Greece, where we studied this species, they were observed to exhibit a facultatively eusocial behavior (Wyman and Richards 2003).

The social plasticity of *Lasioglossum* and its potential as a model for social evolution leads us to believe that observing this group of bees can give invaluable insight on how social behavior affects biological clocks. To test our idea, we captured

these wild bees as they were visiting flowers and observed them in the laboratory using constant conditions. Based on these observations, we determined the variety of behaviors present, and made inferences on how they associate with their sociality.

Methods and Materials:

Study sites

Puerto Rico:

Lasioglossum ferrerii and *L. enatum* were captured using 15 mL falcon tubes from flowers at the Balneario de Luquillo parking lot (18.38706 N 65.72517W, 3 Meters) in Puerto Rico. This site is characterized by having many vine-type plants, high vegetation density, and it is further located right next to a road where *Momordica charantia* is quite abundant. (Figure 1.A and 1.B). Most bees were caught between 8:00 and 12:00h at the flowers of *Momordica charantia* L. (Cucurbitaceae), *Sida acuta* Burm.fil. (Malvaceae), and *Bidens alba* (L.) DC.(Asteraceae). We also observed them visiting *Euphorbia heterophylla*, which has not been reported in previous literature. Collections took place during the months of December, January, March, and August. In total, we collected 36 bees, 26 of which were *L. ferrerii* and 10 were *L. enatum*.

Greece:

Systropha curvicornis and *L. malachurum* bees were collected between 6:00 and 9:00 h from flowers of *Convolvulus arvensis* (Convolvulaceae) that were growing on a recently cut wheat field in Skala Kallonis (39° 10'N 26° 20'E, 0 Meters) on the

548 Island of Lesbos, Greece. We used 15 mL falcon tubes to catch bees in the field,
549 which would house the individual for the duration of the experiment. Sampling was
550 conducted on July 3 of 2017. From this sampling 118 bees were *L. malachurum* and
551 34 were *S. curvicornis*.

552 *Laboratory settings*

553 After collection in the falcon tubes, bees were provided with food that lasted for the
554 whole of the observation period. The food recipe we used varied between the studies
555 in Puerto Rico and Greece (As explained below). The main nutrient for both recipes
556 was sugar, and are therefore nutritionally comparable. However, the agarose based
557 recipe (Puerto Rico) was more convenient in terms of ease and speed of preparation
558 due to the fact that an independent water system was not necessary.

559 Food preparation varied by locations as follows: In Puerto Rico, for every 0.89 ml of
560 water, 1 g of sucrose and 0.1 g of agarose were used. The water was heated in a
561 stirring plate with a magnetic stirrer placed at the bottom. We added the sucrose first
562 to the solution. The agarose was then incorporated upon its dissolution. We left the
563 solution stirring until it turned into a lighter color while being mindful of not letting the
564 solution overheat, as to avoid part of its volume being lost to evaporation. As a form
565 of assurance, we made 3 ml more than what was expected to be used. After all
566 solids were diluted, we quickly pipetted 1 ml into the bottom of a 15 ml centrifuge
567 tube, being mindful of not letting it splash, as to preserve as much solution as
568 possible. Once all of the tubes had their portion of the solution, they were allowed to
569 reach room temperature, and were finally refrigerated. The final product was a gel

570 that could be kept refrigerated until the day it needed to be used as long as it did not
571 dehydrate.

572 In Greece, captured *S. curvicornis* and *L. malachurum* were fed with ApiYem
573 (Namik Kemal University with Kosgeb R&D Innovation Project) which is a food
574 substitute composed of 78.5% sugar and 21.5% invert syrup. Food was placed in the
575 cap-end of each tube and a damp cotton was placed in the other end of the tube to
576 provide water to the bees. The water supply was refilled every 2–3 days. Resources
577 were provided *ad libitum* during the complete running period of the experiment.

578 *Locomotor activity monitoring*

579 Each bee was monitored individually for at least seven days in the falcon tube
580 in which they were captured. The tubes were plugged using cotton balls to let air
581 circulate. These tubes were then placed into TriKinetics' Locomotion Activity
582 monitors (LAM16) that, in turn, were put inside incubators that were set to constant
583 conditions. In Greece, the temperature was 26°C, humidity was at 78%, and light
584 conditions were constant darkness. In Puerto Rico, the temperature was 30°C,
585 humidity was at 65%, and light conditions were constant darkness. The differences
586 in the environmental chamber conditions were set to resemble the average daytime
587 parameters at each location.

588 *Species Identification*

589 The individuals caught in Puerto Rico were identified using Gibbs (2018).
590 Samples collected in Greece were identified by an in-field expert, Victor H. Gonzalez
591 (University of Kansas).

Data processing and analysis

Circadian Analysis

Circadian rhythm and locomotor activity for our subjects were analyzed using the MATLAB toolboxes developed in Jeffrey Hall's laboratory (Levine et al. 2002). The outputs provided data on the individual's locomotor activity throughout the experiment in the form of an actogram, average activity plot, and an autocorrelation that also calculates rhythm strength.

To test if the observed differences in circadian patterns across species, we applied a Brown-Forsythe one way ANOVA with a Dunnett's T3 multiple comparisons test using GraphPad Prism version 8.4.3 for Windows, GraphPad Software, San Diego, California, USA, www.graphpad.com. The variables taken into consideration for this were time, species, individuals, and interspecies variation.

Results:

During the Summer of 2017 on the 3rd of July between the hours of 6:00 am and 9:00 am, *S. curvicornis* and *L. malachurum* were caught as they visited the flowers of *Convolvulus arvensis* on the island of Lesbos. They were transported from the field to the laboratory and placed inside an incubator for 10 days of which 8 were in constant conditions (26 °C and 57% humidity) with the purpose of characterizing their intrinsic biological clock. During this collection, 118 bees were *L. malachurum*, although only 98 survived until the end, and 34 were *S. curvicornis* of which only 4 females survived the study period.

613 Figure 2.A illustrates the average of individuals evaluated under constant
614 conditions of the Greek, solitary and specialist pollinator *S. curvicornis*. Its period
615 runs slightly short at approximately 22 hours, with the peak of its activity in the early
616 morning and an average rhythm strength of 4.4. The population average is
617 consistent with the individuals examined, as illustrated in Figure 2.B; a randomly
618 selected individual looks fairly similar to the activity plotted for the population
619 average. The average for period was 22.75 with a standard deviation of 0.41,
620 Rhythm Strength had an average of 4 and standard deviation of 0.707.

621 The consistency displayed by our population of female *S. curvicornis* is contrasted
622 with the diversity observed in the other 3 species analyzed in this study. This is
623 especially true for the eusocial *L. malachurum*, for whom after careful evaluation of
624 the data we had to create a classification schematic (Figure 3) to appropriately
625 describe the phenotypes being displayed by the population. The population average
626 shows that this species has a perfectly circadian 24-hour period under constant dark
627 conditions. Peak average activity of *L. malachurum* is at 6:00 h, with no clear rest
628 periods when all individuals are averaged. When examined individually, we found
629 five distinct patterns of circadian activity patterns (Figure 3.A and Figure 4). These
630 patterns can be divided into 2 large branches (Figure 3.A), those that are rhythmic
631 and those that are arrhythmic, i.e, individuals with uniform distribution of activity.
632 Rhythmic individuals varied in the amplitude of their activity rhythm, and were
633 therefore classified as strong or weakly rhythmic. Moreover, both strong and weak
634 categories are subdivided into unimodal or bimodal based on the number of activity
635 peaks per day. For example, a bimodal individual is active during two different

instances of the day, such as in the morning and in the afternoon (Figures 4.B and 4.C), while a unimodal individual is mostly active during a set time of the day (Figures 4.D and 4.E).

Strongly rhythmic individuals (either unimodal or bimodal) constituted 41% of individuals. These patterns are recognized by a strong Rhythm Strength (RS)(Figures 4.B.iii and 4.D.iii) on average higher than 2.67 and clear rest and active periods both in the double plotted actogram (Figures 4.B.i and 4.D.i) and the average activity plot(Figures 4.B.ii and 4.D.ii). Weak rhythmicity (Figure 4.C and 4.E) was observed in 21.6% of individuals and they were characterized by having RS values (Figure 4.C.iii and 4.E.iii) that average on 1.79, but their actograms (Figure 4.C.i and 4.E.i) do not show a clear pattern of locomotor activity. Finally, 38% of individuals were arrhythmic. Both the double plotted actogram (Figure 4.F.i) and the average activity plot (Figure 4.F.ii)for arrhythmic bees do not have any discernible daily pattern of activity or inactivity. Often the autocorrelation (Figure 4.F.iii) does not return any values.

On February 19, 2020, between 8:00 am and 10:00 am, at the Balneario de Luquillo (Figure 1.A and B), 36 bees were captured as they visited *Bidens alba*, *Momordica charantia* and *Sida acuta*. Bees were captured and monitored individually in one tube each (modified from Giannoni-Guzman et al. 2014). In the laboratory, the bees were placed inside an incubator for seven days in constant conditions so we could characterize their intrinsic biological clock.

Of the 36 bees captured, 26 were identified as *Lasioglossum ferrerii* (Figure 1.C and 1.D) and ten as *Lasioglossum enatum* (Figure 1.E and 1.F). Only 22 *L. ferrerii* and 8 *L. enatum* survived the entire observation period and were used for analysis. The average peak of circadian activity for *L. ferrerii* is between 6:00–7:00 (Figure 5.A.ii) with a 23-hour period (Figure 5.A. iii), making it rather short. Individuals fell into two categories: 50% were strongly rhythmic and 50% were weakly rhythmic (Figures. 5B and 5C). The peak of average activity for *L. enatum* is from the fifth to the seventh hours of the day, with a circadian period of 23.8 hours (Figure 6.Aiii), just slightly short of one day. The average peak of activity for *L. enatum* was 6:00-7:00 (Figure 6.A.ii). *L. enatum* also had three patterns of activity with similar characteristics to that of the umbrella categories for *L. malachurum*, and we saw it fit to categorize them in a similar fashion (Figure 6). 12.5% of the observed population fell into the strongly rhythmic category, 25% in the weak rhythmic category, and 62.5% in the arrhythmic category.

In summary, all four of the described species followed unique patterns of behavior (Figure 7.A) characterized by the amount of interindividual variation. *Systropha curvicornis* was the species with the least amount of observed interindividual variation in its daily activity patterns, followed by *L. ferrerii* with two distinct patterns of behavior, then *L. enatum* with 3, and lastly, *L. malachurum* with 5.

Cross species comparison of observed circadian parameters.

Average activity was only significant between *L. ferrerii* and *L. malachurum* with a p-value of 0.0016 (Figure 7.B). Circadian Period (Figure 7.C) on the other

679 hand showed differences between *S. curvicornis* and *L. ferrerii* with a p-value of
680 0.0102 as well as *L. malachurum* and *L. ferrerii* with a p-value of <0.0001. Lastly,
681 with even more differences still, Rhythm Strength (Figure 7.D) presented differences
682 between *S. curvicornis* and *L. enatum* (p-value = 0.0014), *S. curvicornis* and *L.*
683 *malachurum* (p-value = 0.0177), *L. ferrerii* and *L. enatum* (p-value = 0.0056) and
684 finally, *L. ferrerii* and *L. malachurum* (p-value = 0.0259).

685 **Discussion:**

686 The solitary specialist, *S. curvicornis* as observed in this work, suggests that
687 at least for the females, the population is consistent, displaying a single circadian
688 activity phenotype (Figure 7.A). The activity of these bees is highly rhythmic, and
689 shows little variation across samples with an average RS of 4.4 and the peak of
690 activity appears to be near the hour 6 of the day. Overall, the species exhibits a short
691 period phenotype under constant darkness. This high degree of rhythmicity might be
692 due to *S. curvicornis*' evolutionary history as a foraging specialist of *C. arvensis*,
693 which blooms for a brief period during the morning, a pattern described for another
694 closely related species (*S. planidens*: Gonzalez et al. 2014). A rigorous internal clock
695 is important to be able to anticipate the time when resources are available. For
696 example, the immediate development of *Osmia bicornis*' circadian rhythm (Beer and
697 Helfrich-Förster 2020) may be related to nourishment accessibility, as it has been
698 shown in the past that large quantities of pollen are the key to proper larvae
699 development rather than diversity of pollen (Radmacher and Strohm 2010). All three
700 of these species (*S. curvicornis*, *S. planidens* and *O. bicornis*) lead solitary lifestyles

and consequently must assure the survival of their progeny in an individual manner. A strong circadian rhythm can ensure that a female may find sufficient resources efficiently to feed its young.

All three species of *Lasioglossum* examined were shown to have more than one distinct pattern of circadian activity. The most diverse of the three was the facultatively eusocial *L. malachurum* (Figure 7.A) with 5 distinct circadian behaviors. Two of the sub-categories of these behaviors fall under the strongly rhythmic category, which we are calling binomial and unimodal. These rhythmic categories are characterized for having an easy-to-distinguish pattern in the actogram, clear rest/activity periods in the average activity plot, and an RS higher than one. Another two of the subcategories fall under the weakly rhythmic umbrella. This umbrella, just like the strongly rhythmic category, can be divided into bimodal and unimodal. These categories can be identified by an actogram with no clear pattern, an average activity plot with more or less clear rest/activity pattern, and an RS larger than one. Lastly, there is the arrhythmic category where no discernible pattern can be pinpointed in the actogram nor in the average activity plot and its RS is less than one. A conceptual map on how these categories are identified can be found in Figure 3.A.

To make descriptions comparable across species, we used the same metrics to categorize the other two bees examined in this study. In categorizing *L. ferrerii* and *L. enatum*, our categories worked as a good basis. *L. ferrerii* only had two distinguishable patterns: rhythmic and noisy rhythmic (Figure 7.A). We decided to change the name from weakly rhythmic to noisy rhythmic because it is a better

723 descriptor (Figure 3.B). Similarly, *L. enatum*, which lives in the same environment as
724 *L. ferrerii*, has 3 distinguishable categories (Figure 7.A). These categories are
725 rhythmic, noisy, rhythmic and arrhythmic. In contrast to *L. ferrerii*, *L. enatum*
726 expressed 5 individuals in the arrhythmic category. Taking into consideration that
727 both of these species of bees were caught in the same environment, and that they
728 belong to the same genus, the results suggest that something other than
729 environmental variables are behind these differences.

730 The difference in expression of circadian patterns between *L. ferrerii* and *L.*
731 *enatum* could be explained by competition. Both of these species share the same
732 niche in Luquillo, to the point of being caught in the same flowers during the same
733 range of time. Having a slight difference in rhythmicity can lower the possibility of
734 temporal competition when foraging. *L. ferrerii* on average would be active from 5:00
735 am to 10:00 am, while average time of activity for *L. enatum* would be from 3:00 am
736 to 8:00 am. Due to that two-hour disphase, it would appear to be less likely that bees
737 from these two species try to visit a flower simultaneously, yet their schedules still
738 have some overlap. These observations are echoed by another study conducted in
739 Greece where they demonstrated that 3 species of carpenter bees (*Xylocopa spp.*)
740 that share the same resources have different circadian rhythms when measured
741 under natural field conditions and also in artificial constant and oscillating conditions
742 (Ortiz-Alvarado et al. in rev.). While the solitary, *Xylocopa* species have interspecies
743 variation in their circadian rhythms, two out of the three examined in Ortiz-Alvarado
744 et al. follow a similar pattern as *S. curvicornis*, where there isn't much, if any
745 individual differences in the populations examined. Therefore, in that particular case,

competitor effects can explain the differences in rhythm across species, but in the case of *L. ferrerii* and *L. enatum*, it cannot explain the individual differences observed at the species level.

At a higher level looking at the statistical analysis of all four halictid bees (Figure 7), some interesting patterns can be noted. In terms of average activity (Figure 7.B), there was only a difference between *L. ferrerii* and *L. malachurum*, none of the other possible combinations of differences occurred. However, the length of the whiskers in the box plots for both *L. enatum* and *L. malachurum* does suggest a level of diversity at the intraspecies level and could be reflective of the number of circadian behaviors observed in these species.

When analyzing the circadian period, the observed differences were between *L. ferrerii* and *L. malachurum* as well as *S. curvicornis* and *L. malachurum* (Figure 7.C). The latter of these pairs have shared environmental conditions when the former pair does not. It is also interesting to note that *L. ferrerii* and *S. curvicornis* cannot be found in the same locations, and yet, they do not appear to have significantly different circadian periods. In fact, for the populations examined, they appear to be comparable.

The picture becomes clearer still when observing the differences in rhythm strength (Figure 7.D). Where those species with a lesser number of circadian phenotypes are more similar to each other, and likewise, those with the most diversity are more similar to each other. In other words, *S. curvicornis* and *L. ferrerii* were both significantly different to *L. enatum* and *L. malachurum*, but not to each

other. Likewise, there was no significant difference between *L. enatum* and *L. malachurum*. Because there is this consistency of differences that is not associated with differences in environments, we believe that the key to explaining the difference in diversity of behaviors may not lay in competition, but in something more endogenous of the species. Nevertheless, more data is needed.

As we mentioned before, *Lasioglossum* as a genus is well-known for having a large diversity in social behaviors. This diversity in sociality may also be reflected in other types of behaviors, and could be the key to explaining the individual differences in circadian rhythm we observed at the interspecies level. A limitation of our work was the sample size and the number of the evaluated populations (one population for each species and low number of individuals, particularly for *L. enatum*), thus, it definitely merits repetition of the work to see if our results can be replicated. Furthermore, in the one species whose population was close to a hundred, more time than usual was needed to evaluate the data, due to the diversity found within it. In future studies we will focus on streamlining the process of describing a species with diverse circadian behaviors such that it will facilitate studies with a higher volume of observations. Additionally, we will continue describing the circadian rhythm of more species of *Lasioglossum* who present social behaviors not evaluated in this study, and will evaluate if there is a relationship between sociality and rhythm, causal or otherwise.

790 **Figure Legends:**

791 **Figure 1: Habitat (A-B) and Species Observed (C-F).** A) Puerto Rico study site in
792 which *L. ferrerii* and *L. enatum* were captured. B) Some of the vegetation the bees
793 were observed visiting, with flowers belonging to the families: Commelinaceae,
794 Cucurbitaceae and Euphorbiaceae. C) Female of *L. ferrerii* distinguished from the
795 male by its short antenna and pointed abdomen. D) Male of *L. ferrerii*, distinguished
796 by its long antennae and flat abdomen. This species is known for its long head
797 shape and metallic metasoma (Gibbs 2018). E) Female of *L. enatum* distinguished
798 from the male by its short antenna and pointed abdomen. F) Male of *L. enatum*,
799 distinguished by its long antennae and flat abdomen. This species is distinguished
800 by: “tegula punctate, extended posteriorly to form a small angle, mesepisternum
801 punctate and metasoma brown” (Gibbs 2018).

802 **Figure 2: Female *S. curvicornis* exhibit short period phenotype under constant**
803 **darkness (<24 h endogenous circadian rhythm).** (i) Double-plotted actogram
804 showing the locomotor activity pattern of: (A) The average of all 4 individuals
805 examined of *S. curvicornis*. (B) A representative individual randomly selected from
806 the population. In a double plotted actogram, each row represents locomotor activity
807 (counts per 30 min) of two consecutive days and the second is repeated such that it
808 is always the beginning of the next row. The x-axis shows the time of day under
809 constant darkness expressed as circadian time (CT). (ii) Average of the locomotor
810 activity patterns of the five days observed into one. The y-axis represents activity
811 over time and the x-axis represents CT (iii) Autocorrelation plots used to determine

812 the period (p), rhythm index (RI) and rhythm strength (RS). The oscillations indicate
813 periodicity. The asterisk at the third peak of the autocorrelation plot indicates the
814 particular time point used for the determination of the rhythm parameters.

815 **Figure 3: A summary of the variations in the circadian rhythm as observed in:**

816 A) *Lasioglossum malachurum*, B) *Lasioglossum enatum* and *Lasioglossum ferreirii*.

817 The circle represents the root of the flowchart, squares represent nodes that branch
818 off and rhombuses represent leaves. In total for malachurum, 5 distinct behaviors
819 were observed.

820 **Figure 4: *L. malachurum* exhibits a variety of circadian phenotypes under**

821 **constant dark conditions.** (i) Double-plotted actogram showing the locomotor
822 activity pattern for 5 days of: (A) The average of all 98 individuals examined and
823 representatives for the following categories: (B) Bimodal Rhythmic (C) Weakly
824 Rhythmic Bimodal, (D) Unimodal Rhythmic, (E) Weakly Rhythmic Unimodal, (F) and
825 Arrhythmic circadian behaviors. (ii) An average activity plot for the five days of
826 observation (iii) Autocorrelation plots used to determine the period (p), rhythm index
827 (RI) and rhythm strength (RS).

828 **Figure 5: Description of the circadian behaviors exhibited by *L. ferreirii* under**

829 **constant dark conditions.** (i) Double-plotted actogram of the locomotor activity
830 from the five-day observational period for: A) All 22 individuals from the data set
831 averaged out into one representative individual. B) A representative individual out of
832 the 11 from the category Strongly Rhythmic. C) A representative individual out of the
833 11 from the category Noisy Rhythmic. (ii) An average activity plot for the five days of

observation (iii) Autocorrelation plots used to determine the period (p), rhythm index (RI) and rhythm strength (RS).

Figure 6: Description of the circadian behaviors exhibited by *L. enatum* under constant dark conditions. (i) Double-plotted actogram of the locomotor activity from the five-day observational period for: A) All 8 individuals from the data set averaged out into one representative individual. B) The only individual from the category Strongly Rhythmic. C) A representative individual out of the 2 from the category Noisy Rhythmic. D) A representative individual out of the 5 from the Arrhythmic category. (ii) An average activity plot for the five days of observation (iii) Autocorrelation plots used to determine the period (p), rhythm index (RI) and rhythm strength (RS).

Figure 7: Summary of descriptive and inferential statistics. A) Number of circadian categories observed by species. B) Box plot illustrating the difference in average locomotor activity between species. *S. curvicornis* has a minimum of 8.200, 25% percentile of 8.900, mean of 13.63, 75% percentile of 20.28 and a maximum of 22.90. *L. ferrerii* has a minimum of 0.7000, 25% percentile of 2.675, mean of 4.950, 75% percentile of 7.100 and a maximum of 11.20. *L. enatum* has a minimum of 1.000, 25% percentile of 2.875, mean of 8.113, 75% percentile of 6.650 and a maximum of 36.00. *L. malachurum* has a minimum of 6.700, 25% percentile of 10.35, mean of 8.201, 75% percentile of 10.35 and a maximum of 26.30. There was only a statistical difference between *L. ferrerii* and *L. malachurum* with a p-value of 0.0016, DF of 61.66 and t of 3.867. C) Box plot illustrating the difference in circadian

856 period between species. *S. curvicornis* has a minimum of 22.20, 25% percentile of
 857 22.35, mean of 22.75, 75% percentile of 23.10 and a maximum of 23.20. *L. ferrerii*
 858 has a minimum of 21.80, 25% percentile of 21.95, mean of 22.69, 75% percentile of
 859 23.20 and a maximum of 24.200. *L. enatum* has a minimum of 20.00, 25% percentile
 860 of 22.40, mean of 23.31, 75% percentile of 24.60 and a maximum of 25.50. *L.*
 861 *malachurum* has a minimum of 20.20, 25% percentile of 23.50, mean of 24.00, 75%
 862 percentile of 24.50 and a maximum of 27.80. Both *S. curvicornis* and *L. ferrerii* were
 863 significantly different from *L. malachurum* with p-values of; 0.0102 and <0.0001, DFs
 864 of; 5.652 and 52.94 and, t of; 5.179 and 6.565, respectively. D) Box plot illustrating
 865 rhythm strength among species. *S. curvicornis* has a minimum of 3.000, 25%
 866 percentile of 3.250, mean of 4.000, 75% percentile of 4.500 and a maximum of
 867 4.500. *L. ferrerii* has a minimum of 0.6000, 25% percentile of 1.700, mean of 2.691,
 868 75% percentile of 4.125 and a maximum of 4.600. *L. enatum* has a minimum of -
 869 2.500, 25% percentile -0.4250, mean of 0.2125, 75% percentile of 1.200 and a
 870 maximum of 2.200. *L. malachurum* has a minimum of -2.100, 25% percentile of
 871 0.7250, mean of 1.754, 75% percentile of 2.800 and a maximum of 4.500. The
 872 solitary, *S. curivcornis*, and communal *L. ferrerii*, were significantly different from the
 873 eusocial species, but not each other. Likewise, *L. enatum* and *L. malachurum* were
 874 not significantly different. *S. curvicornis* vs. *L. enatum* ; p-value of 0.0014, df of 7.932
 875 and t of 6.226. *S. curvicornis* vs. *L. malachurum*; p-value of 0.0177, df 4.04 of and t
 876 of 5.895. *L. ferrerii* vs. *L. enatum*; p-value of 0.0056, df of 17.09 and t of 3.982. *L.*
 877 *ferrerii* vs. *L. malachurum*; p-value of 0.0259, df of 33.42 and t of 3.054.

878

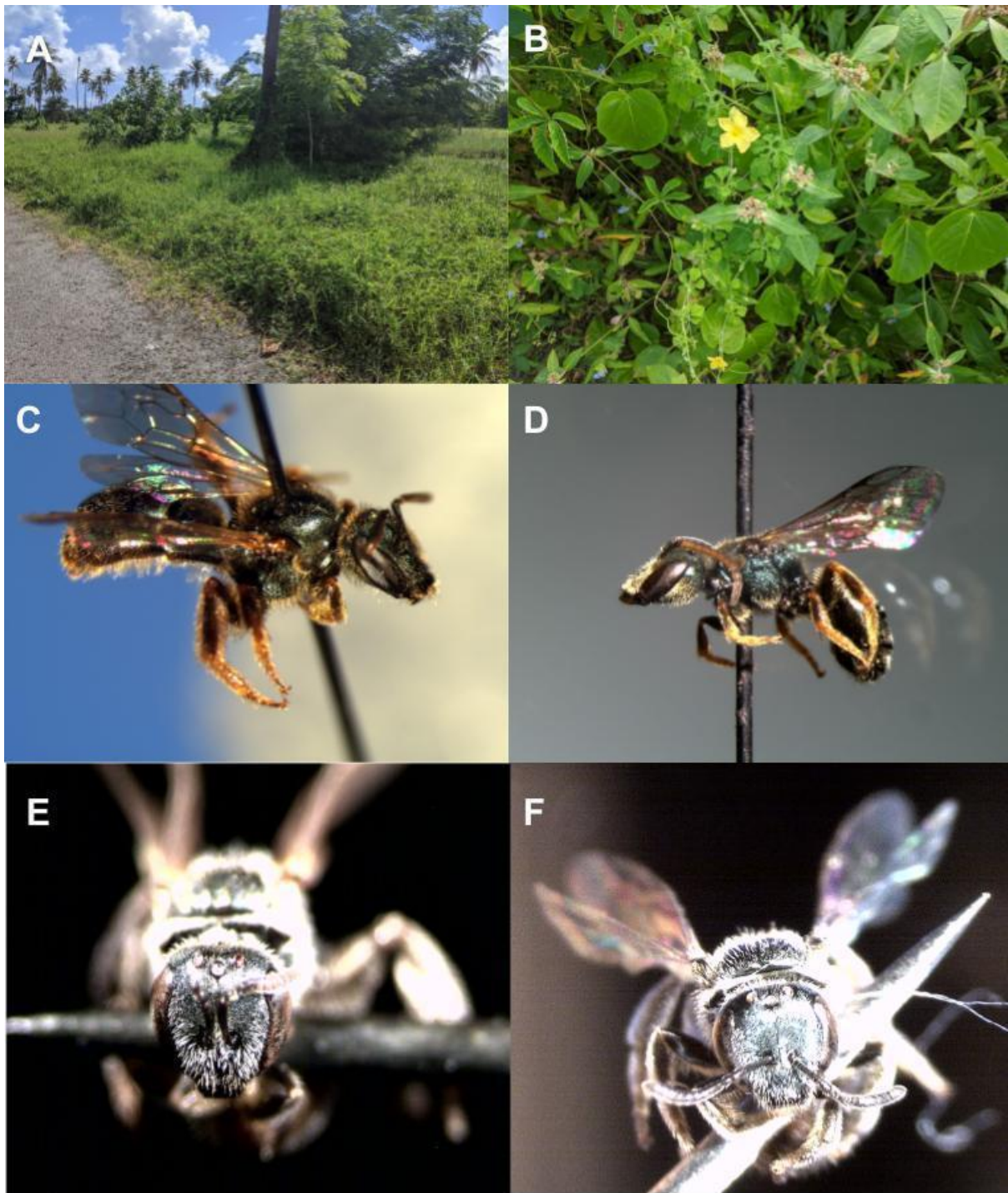


Figure 1: Habitat (A-B) and Species Observed (C-F).

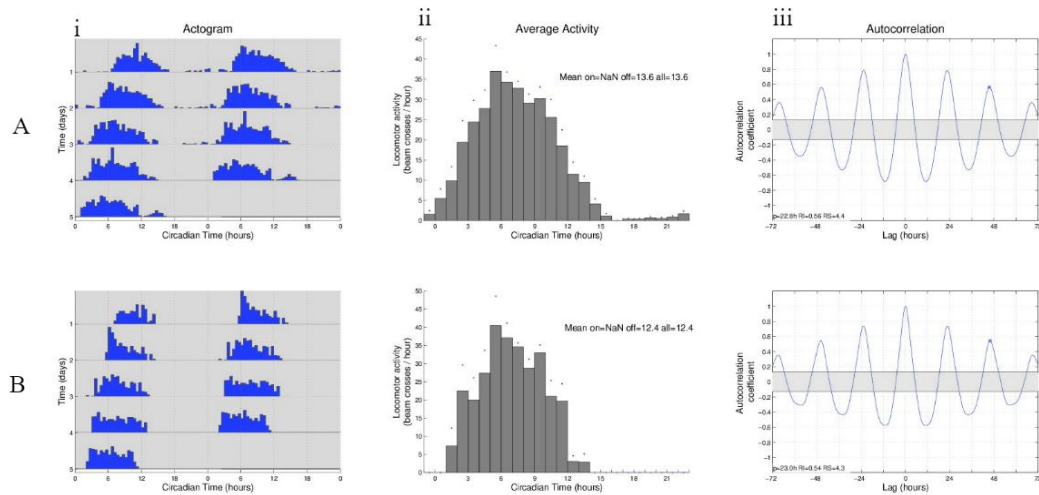


Figure 2: Female *S. curvicornis* exhibit short period phenotype under constant darkness (<24 h endogenous circadian rhythm)

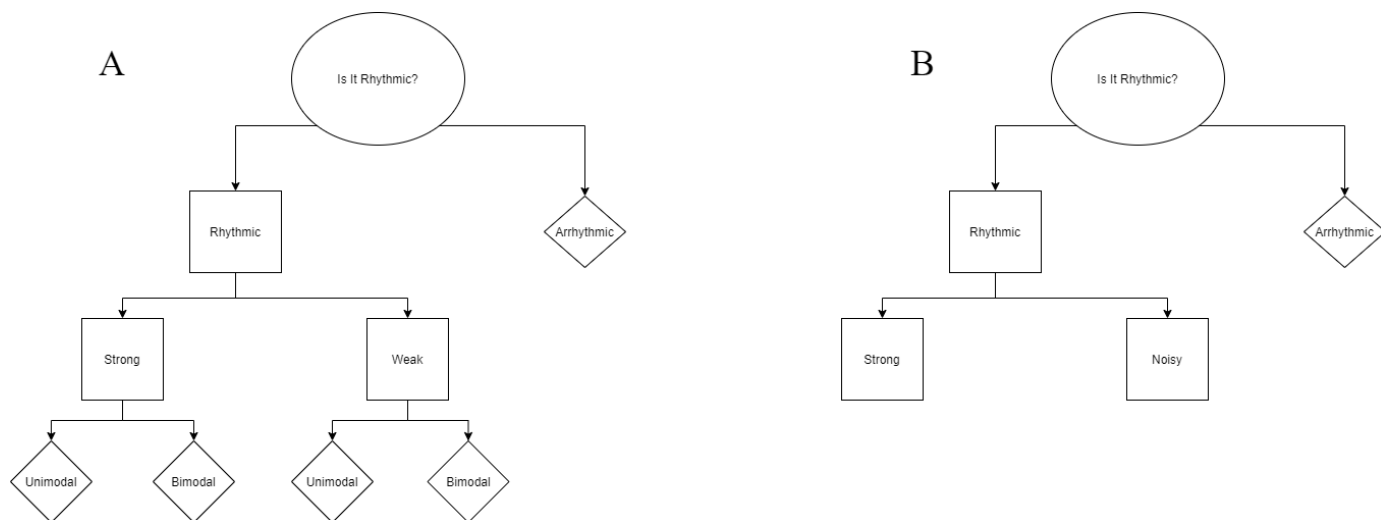


Figure 3: A summary of the variations in the circadian rhythm as observed in:
 A) *Lasioglossum malachurum*, B) *Lasioglossum enatum* and *Lasioglossum ferreirii*

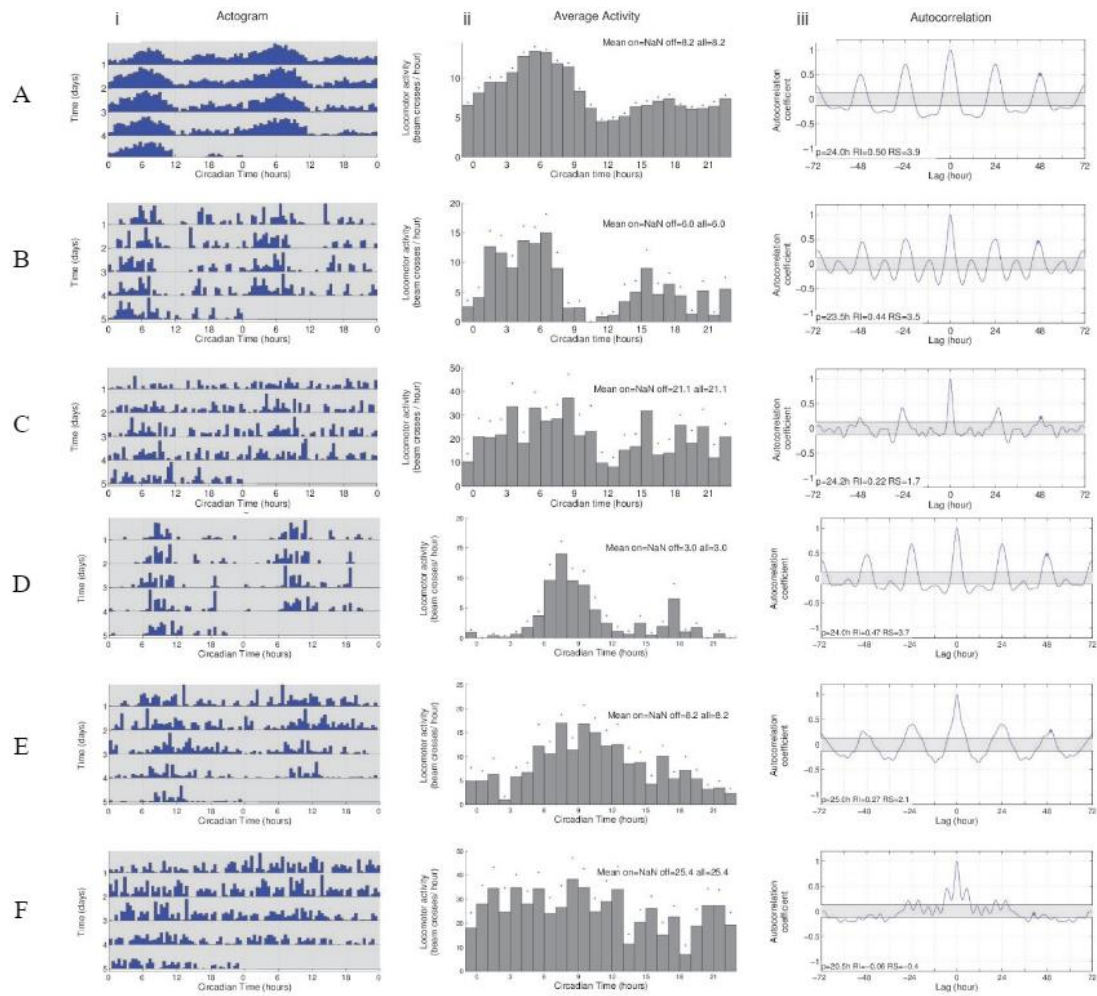


Figure 4: *L. malachurum* exhibits a variety of circadian phenotypes under constant dark conditions.

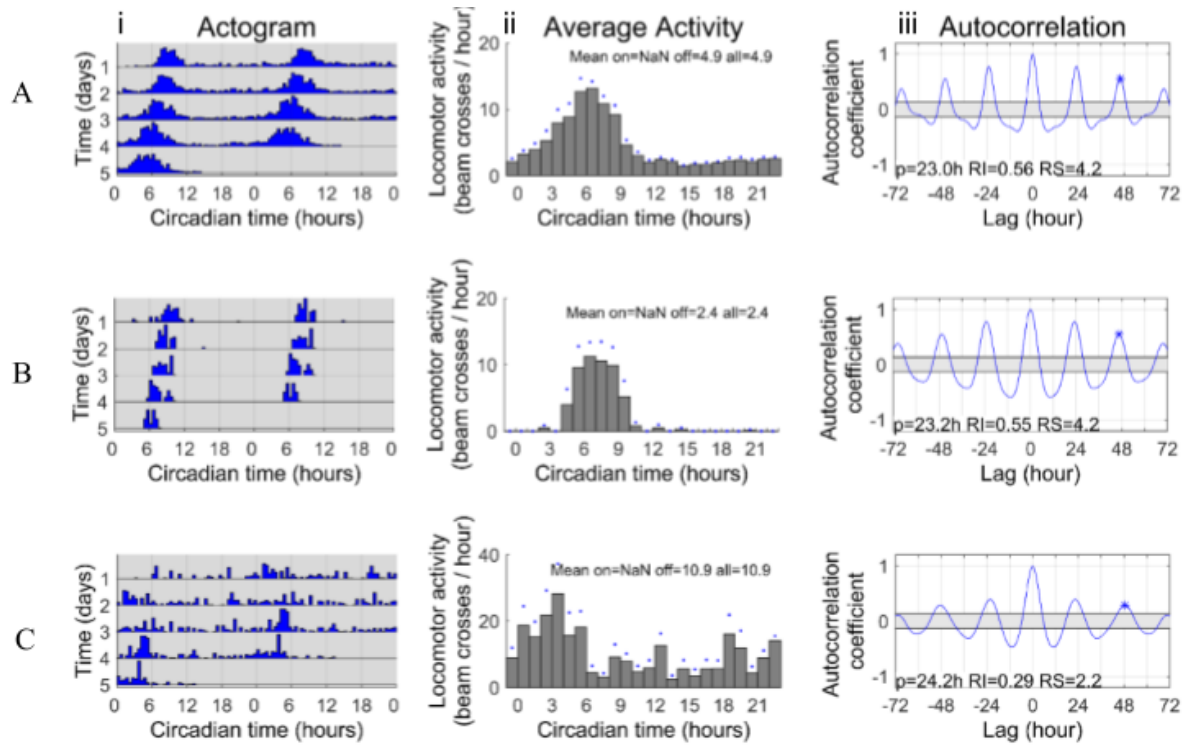


Figure 5: Description of the circadian behaviors exhibited by *L. ferreirii* under constant dark conditions

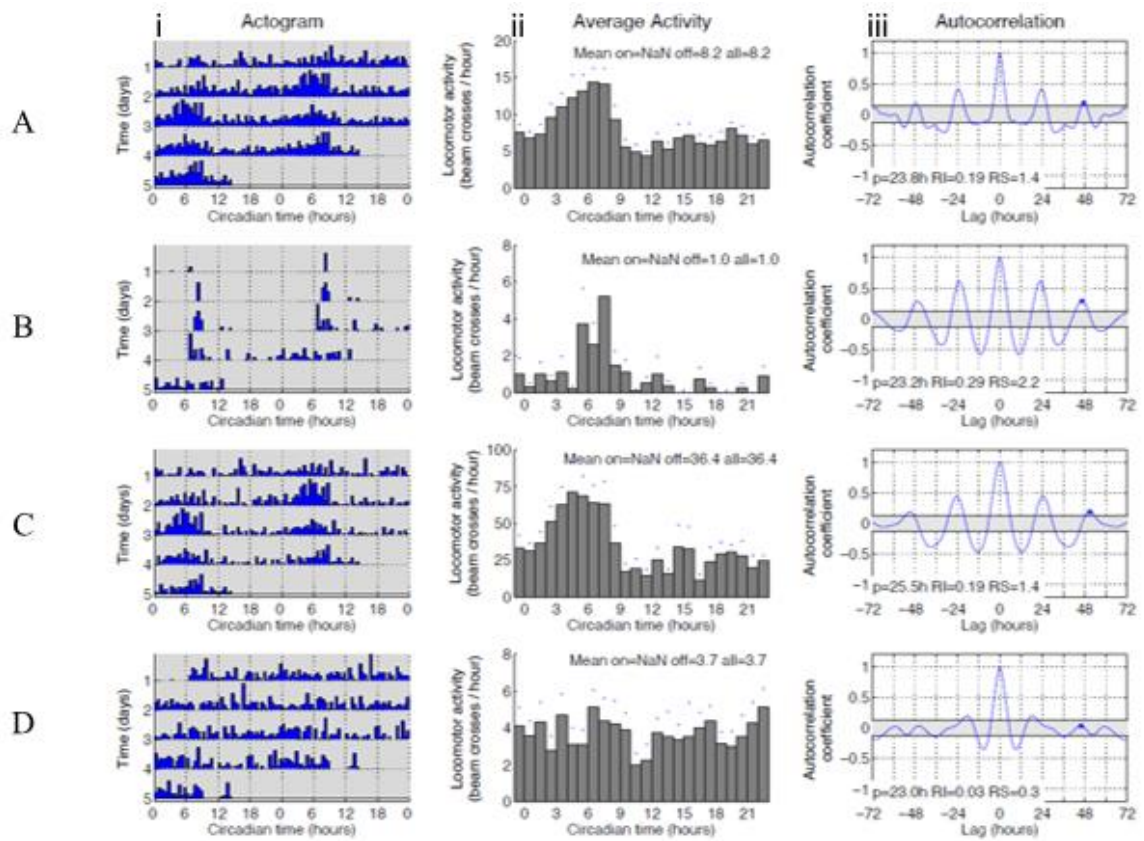


Figure 6: Description of the circadian behaviors exhibited by *L. enatum* under constant dark conditions.

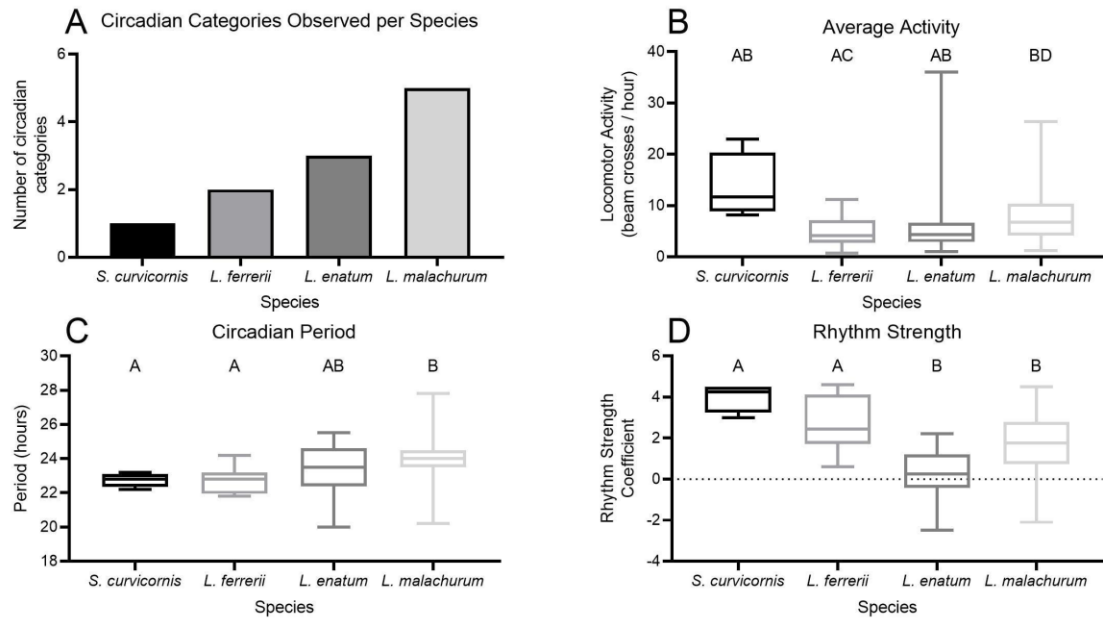


Figure 7: Summary of descriptive and inferential statistics

916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940

Chapter 2

A Machine Learning Pipeline for the Classification of Inter-individual Behaviors in Circadian Rhythms of *Lasioglossum* Bees

A Machine Learning Pipeline for the Classification of Inter-individual Behaviors in Circadian Rhythms of *Lasioglossum* Bees

Abstract:

Traditionally, model organisms are used for the study of daily activity patterns. Because the genetics of these organisms are well-known, there has been no need to do systematic *a priori* sorting of the individuals into groups, the categories are built into the analysis pipeline. Recently, there has been a rising interest in using non model organisms in studies, and with them comes unexplored diversity. To facilitate the *a priori* sorting of individuals for analysis, we designed a machine learning pipeline using individuals of *Lasioglossum malachurum* as a case study. We tested both supervised and unsupervised algorithms, and evaluated how well they separated the individuals of the population in discrete groups based on the phenotype of their daily activity patterns. Decision Trees with Symbolic Aggregate approXimation (SAX) transformation achieved the best results. K Nearest Neighbors with Symbolic Aggregate approXimation was a close second. In the future, we aim to test the transferability of this pipeline using insects of the same genus, but different expressions of behavior.

Introduction:

The study of circadian rhythms focuses on characterizing the daily activity patterns of populations. These daily activity patterns are constrained to a period of time that approximates 24 hours and are consistent with the rotation of Earth. Traditionally, model systems, such as fruit flies, mice, rats, and *E. coli* cell lines, are used to evaluate circadian rhythms. Because these models are well-known in terms

of phenotype, no *a priori* systematic sorting of the data sets is usually done, the scientist already knows which individuals have which genotypes and can sort them accordingly.

Recent works in behavioral biology and circadian rhythms have evaluated the use of non-model systems, such as varying species of bees and have found them to be ideal systems to study due to the diversity of behaviors exhibited (Bloch and Grozinger 2011). In particular, halictid bees have been suggested for studies that evaluate the impact of sociality on other behaviors, for example circadian rhythms (Toth and Rehan 2016, Bloch and Grozinger 2011, Danforth et al. 2003). However, because bees express a high diversity of behaviors, they present an added layer of difficulty in terms of systematic evaluation. For example, Gianoni-Guzman and collaborators (2014) saw social insects in comparison to *Drosophila melanogaster* (the fruit fly) displayed a significant amount of variance in terms of period. This means that the length of the internal clock of these social insects (which includes honey bees) varied from individual to individual. In a different work, Gianoni-Guzman (2020) tracked the foraging schedules of honey bees. He found three distinct temporal shift behaviors. The variance in the first study, where one of the populations evaluated was of honey bee foragers, can be explained by the behaviors observed in the second. Studies like Gianoni-Guzman 2014, where there is no prior knowledge of the phenotypes present in a population being evaluated, may benefit from *a priori* sorting based on similarity found within subsets of the population. Studies like that of Gianoni-Guzman (2020) may be of use to train a classification pipeline and serve as a sorting tool for future studies. In cases where no prior labels

are available, a clustering pipeline in conjunction with expert validation can inform the labeling process. Expert validation is defined as having an expert evaluate the clusters and validate their biological significance.

In this work, we evaluate the addition of a grouping step to the circadian evaluation pipeline, as to facilitate the determination of the groups in the non-model systems. To achieve this, we applied both clustering (K-mean and Partitioning Around Medoids) and classification (K Nearest Neighbors, Decision Trees and Random Forest) on a locomotion activity dataset for the halictid bee *Lasioglossum malachurum* (*L. malachurum*). In addition, we transformed the time series of each individual using Symbolic Aggregate approXimation (SAX) to reduce dimensionality in the data set, so that the computational resources needed to evaluate this data set are minimized. Given how circadian activity is traditionally analyzed, we expect that the methods that best capture the shape of daily activity patterns will be the most effective at grouping the bees in groups of circadian significance.

High Level View of the data:

The following 3 sets of graph panels in Figure 1 are examples of what the data looks like after it's analyzed with a traditional circadian pipeline. Visually, their activity patterns do not look similar. The individual in Figure 1. A is extremely rhythmic, which makes it so the graph of the average daily activity looks fairly similar to the activity of any one day in the actogram. This individual appears to have a bimodal pattern of activity, where the main bout of activity occurs in between the sixth and twelfth hours of the day, and the minor one occurs about an hour later after a period of inactivity. The Lomb Scargle Periodogram suggests that a day for this

1010 individual last about 23.77 hours, and the autocorrelation oscillates, which indicates
1011 strong rhythmicity.

1012 In contrast, the individual in Figure 1. B appears to be a repeated yet spread
1013 out bimodal pattern. That spreading of the locomotor activity already makes it
1014 different from the individual in Figure 1. A. They also differ in the time of day in which
1015 they are active. The average activity plot in Figure 1. B has a similar shape to the
1016 days of observation, but there is a distinct difference between it and Figure 1.A.. One
1017 can see that the gap in activity isn't exactly the same for every day. The
1018 Periodogram suggests that this individual's day lasts 22.85 hours, which is slightly
1019 shorter than the previous individual's day. The autocorrelation graph does not have
1020 as well-defined undulations as the previous individual, which suggests that the
1021 individual Figure 1. B, although rhythmic, is not as strongly rhythmic as the subject in
1022 Figure 1. A.

1023 Lastly, Figure 1. C is an example of an individual whose rhythm is not
1024 circadian (24h), but might be ultradian (12h). Different from the other two Figures, it's
1025 difficult to determine a pattern from looking at the actogram, and by consequence,
1026 it's not easy to tell if the average is in any way representative of the individual. This
1027 makes us rely on the autocorrelation to determine rhythm. As one can see in the
1028 autocorrelation plot, this one has a smaller undulations pattern, although it rarely
1029 reaches significance. Those small undulations can lead one to believe that this
1030 individual might have a much shorter internal day than what we are accustomed to
1031 analyzing. Therefore, the periodogram might not be useful to determine the period
1032 for these cases.

In all, these three individuals demonstrate that there exists a level of diversity within the data set that might be lost if we evaluate it only using averages. In Figure 2, this theory is illustrated with a single plotted actogram, which represents the activity of all 98 individuals present in this data set. The activity of any one of the days in Figure 2 is not representative of any one day illustrated in the individuals in Figure 1.

To report averages as representative of the species in terms of circadian rhythm, in this case, is misleading. Nevertheless, describing each individual is tedious work. As an alternative, we should be able to group together individuals with similar characteristics and use them to describe the daily activity patterns of the species. To achieve this, we explored the use of both unsupervised and supervised machine learning with an abstract knowledge representation of time series. Because the methods used to evaluate circadian data are dependent on the mean and shape of the activity, we expect that methods that utilize centroids will be the most effective, and we chose SAX for its ability to cluster and classify univariate time series by shape.

Methods:

Animal Model:

Lasioglossum malachurum (*L. malachurum*) (Kirby, 1802)

L. malachurum is an obligately eusocial halictid bee, also known as sweat bees. They typically nest underground in complex colonies (Wyman, L. M., & Richards, M. H. 2003).

Data collection:

1056 *Site description:*

1057 The specimens were collected on Lesbos (39°10'N 26°20'E), a Greek island
1058 in the northern Aegean Sea off the coast of Turkey. It was summer during the time of
1059 collection, and the specimens were captured after being observed between the
1060 hours of 0600 and 0900 (Cordero-Martínez, C.S., et al. 2017).

1061 *Capture methods*

1062 The bees were hand captured as they visited the flowers of *Convolvulus*
1063 *arvensis*, colloquially known as morning glories (Cordero-Martínez, C.S., et al. 2017).

1064 *Housing and observation*

1065 Each bee was housed individually in a modified centrifuge tube, first in
1066 oscillating conditions and later in constant conditions. Oscillating conditions were
1067 meant to mimic the changes of light, temperature and humidity in their natural
1068 environment, while constant conditions kept all environmental cues constant. While
1069 in constant conditions, the bees were not exposed to light. At the bottom of the tube
1070 that housed the bee, there was a cotton ball soaked in water for hydration that was
1071 refilled every 2-3 days. The body of the tube had small holes to allow for air
1072 circulation. On the cap, there was a paste that would function as food for the bee.
1073 The paste was composed of a modified version of ApiYen brand bee feed that had
1074 no protein, but kept the same amount of sugar.

1075 These tubes were placed into TriKinetics' Locomotion Activity monitors (LAM),
1076 which in turn were put inside GRW-20 CMP3/TBLIN incubators. The incubators were
1077 set to mimic the environmental conditions in which *L. malachurum* was captured
1078 (Cordero-Martínez, C.S., et al. 2017).

1079 **Data exploration:**

1080 The code for this process is available at:

1081 <https://github.com/ComplejoC/CircadianThesis>

1082 *Data formatting*

1083 The data is outputted by the LAM as a data frame contained within a tabular
1084 separated value file (.tsv), as specified by the manufacturer in the user manual
1085 (<https://www.trikinetics.com/Downloads/DAMSystem3%20Software%20Data%20Sheet.pdf>). In this data sheet, columns 11-42 are representative of one individual, and
1086 each row is the number of times that the subject moved, as detected by the sensor in
1087 one minute.

1089 The LAM system is prone to a number of errors, of which we must account for
1090 while analyzing the data. If a subject were to fall asleep or otherwise become
1091 immobile directly on top of the beam, the system may display the individual as more
1092 active than they are in actuality. If someone were to open the doors of the incubator,
1093 the sudden entrance of an outside light source can break the beam, and it will count
1094 as activity. In our case, we do not need to worry about the lights inside the incubator
1095 due to the observations being done in Dark-Dark conditions (DD). Nevertheless, it's
1096 important to keep in mind that turning on the lights inside the incubator may cause
1097 false activity counts.

1098 *Data processing*

1099 To be able to properly represent the data, we had to remove dead individuals.
1100 Including them in the study would influence the shape of the data, and might not
1101 accurately represent the actual behavior of the individual. We used the death

1102 detection algorithm from the Rethomics framework of R packages (Geissmann Q,
1103 Garcia Rodriguez L, Beckwith EJ, Gilestro GF. 2019).

1104 When graphing autocorrelations, it is common to detrend to reduce the
1105 influence of distortion that is inherent in this type of data. To do this, we used the
1106 detrend function from the pracma R package (Hans W. Borchers 2019). Because
1107 we wanted to see how the data correlated with itself in both directions, we used the
1108 cross correlation function from base R (ccf)(R Core Team 2019) and plotted it using
1109 the autoplot function from ggfortify (Yuan Tang 2016), resulting in a bidirectional
1110 autocorrelation instead of the unidirectional autocorrelation available in base R.

1111 *Visualization:*

1112 The double plotted actogram, and the Lomb–Scargle Periodograms were all
1113 generated using the Rethomics pipeline (Geissmann Q, Garcia Rodriguez L,
1114 Beckwith EJ, Gilestro GF. 2019).

1115 **Double Plotted Actogram:**

1116 This type of visualization is typically used to represent the Circadian
1117 locomotor activity rhythms. Each vertical bar is representative of
1118 activity, in our case the number of times the bee broke the sensor laser.
1119 The higher the bar, the more active the individual. It's called double
1120 plotted, because with the exception of the first day (and maybe the
1121 last), every day is plotted twice, two consecutive days are plotted one
1122 next to each other, and the second day being re-plotted in the right half
1123 right under (Jud, C. et al., 2005).

1124 **Lomb-Scargle Periodograms:**

1125 Lomb-Scargle is an algorithm that can be used to describe the period
1126 of unevenly sampled time series data sets. It allows for approximation
1127 of a power spectrum estimator, similar to that of a Fourier transform.
1128 The resulting estimators can be used to determine the period of
1129 oscillation of a data set (Jacob T. VanderPlas, 2018). In our case, those
1130 periods are representative of the length of the day in the internal clocks
1131 of our subjects. For ease of visualization, the data was downsampled
1132 from 1 minute bins to 30 minute bins. This down-sampled data was
1133 used to make average activity plots using ggplot2 (H. Wickman. 2016)
1134 and the autocorrelation plots.

1135 **Average Activity plots:**

1136 These plots are representative of the average activity done by each
1137 individual for the duration of the study. We added all days together into
1138 one representative day, and divided it by the total time accumulated to
1139 get the average.

1140 **Autocorrelation plots:**

1141 The autocorrelation of a time-series measures how similar a time-
1142 series is with a forward or backward shifted version of itself. For signals
1143 that oscillate perfectly, the graph of this function oscillates between +1
1144 and -1, $t=0$ is the highest value registered. To create our
1145 autocorrelation plots, we detrended our data using pracma (Hans W.
1146 Borchers 2019) and did a self-cross correlation using the stats ccf

1147 function (R Core Team 2019) to see how the data matched itself in both
1148 directions.

1149 **Clustering Methodology:**

1150 Clustering is a type of unsupervised machine learning, in which the data
1151 analyzed is unlabeled. That is to say, the machine does not have a reference of how
1152 the data groups together, and must learn to do so without input from a user (Géron,
1153 2018). The purpose of this study is to optimize the already existing pipeline for
1154 circadian analysis. Therefore, the data transformations used to reduce
1155 dimensionality for clustering are those that are already used in circadian science:
1156 The Lomb Scargle (LS) Periodogram, Average Daily Activity and Autocorrelations.

1157 *Consensus Clustering*

1158 Before attempting to cluster, we did a procedure known as consensus
1159 clustering. It is a “method to represent the consensus across multiple runs of a
1160 clustering algorithm to assess the stability of the discovered clusters” (Monti, S.,
1161 2003). The use of the word stability in these instances refers to how much the
1162 composition of each cluster changes over repetitions.

1163 We wanted to input a number of clusters for our unsupervised algorithm that
1164 were not arbitrary. By doing various repetitions of the clustering procedure and
1165 observing which number of clusters has the minimum error, we are allowing the data
1166 to speak for itself instead of using our prejudice. For this, we used Wilkerson and
1167 Niel’s R package consensusclusterplus that can be found in bioconductor
1168 (Wilkerson, D. M and Hayes, Neil D 2010). We used a maximum k of 10 and 10,000

1169 repetitions, using both Euclidean and Manhattan distances. Ultimately, the
1170 Manhattan distance was used to minimize cluster collisions.

1171 *K-means clustering*

1172 In one instance, we use K-means to cluster our data. The algorithm takes an
1173 n number of observations or individuals in our case, and divides it into K clusters. K
1174 is the number of clusters as provided by the user. The original centroids around
1175 which the data clusters are randomly set and require a seed for replicability in the
1176 code. We used the set.seed function from Base R (R Core Team 2019) to generate
1177 the centroids, and ran the kmeans function from the stats package (R Core Team
1178 2019). Once the first centroids are set, each individual becomes a member of the
1179 cluster with the nearest mean, serving as a prototype of the cluster. Then the
1180 algorithm continues, by adjusting the centroids until a partition of the data that
1181 minimizes the sum of squares deviation is found (H.-P. Kriegel et al. 2017).

1182 *PAM (K-Medoids)*

1183 In two instances, to cluster our data, we used the PAM algorithm, also known
1184 as k-Medoids, from the R package cluster (Maechler, M et al., 2019). This algorithm
1185 is considered a more robust version of K-means, because it minimizes a sum of
1186 dissimilarities instead of a sum of squared euclidean distances. Just like in K-means,
1187 the algorithm is given a K number of centers, the difference being that centers or
1188 medoids are actual points within the data set (L. Kaufman P. J. Rousseeuw, 1990).

1189 We clustered the data using 3 different transformations: Average Daily
1190 Activity, the AutoCorrelation Coefficient and The Lomb–Scargle periodogram.

1191 *Cluster visualization*

1192 Because we have 97 individuals with 5 days' worth of locomotion data,
1193 graphing the clusters means we have to reduce dimensions to make them human
1194 readable. To do this, we used the R package factoextra (Alboukadel Kassambara
1195 and Fabian Mundt, 2020), which graphs data by using the two most representative
1196 dimensions of the whole and using them as the x/y-axis.

1197 **Classification Pipeline:**

1198 **Phase 1: Creating the Gold Standard Part 1 Building the Baseline**

1199 1. Naive Classification:

- 1200 a. Teach a group of non-experts how to interpret the graphs usually used
- 1201 for circadian analysis.
- 1202 b. Separate them into smaller groups and give them a stack of graphs
- 1203 that they will divide into 3 groups: Rhythmic(R), Arrhythmic (AR) and
- 1204 Weakly Rhythmic (WR).
- 1205 c. Measure the percentage of classification coincidence between the
- 1206 groups.

1207 2. Experienced (Round two using the same people):

- 1208 a. Give a review on how to classify the graphs.
- 1209 b. Separate them once again into smaller groups and give them the same
- 1210 stack of graphs that they will divide into the same 3 groups as last time.
- 1211 c. Measure the percentage of classification coincidence between the
- 1212 groups.
- 1213 d. Compare the percentages of coincidence of the groups between
- 1214 rounds.

1215 **Phase 3: Data Transformation**

1216 Time series data is famously noisy and unevenly sampled. To remedy the
1217 pitfall of using this type of data, one may use transformation methods to make it more
1218 manageable, such as Fourier Transformations, Discrete Wavelets, and Lomb Scargle
1219 periodograms (Refienetti et al. 2007). We decided to use Symbolic Aggregate
1220 ApproXimation or SAX for short. Our choice was based on the capacity of this
1221 algorithm to do dimensionality reduction while still staying true to the data. This goes
1222 in hand with the other advantage of using SAX, which is that it has lower bounds (J.
1223 Lin, 2007). This means that it has the capacity to represent our time series while
1224 using a minimum amount of resources.

1225 1. **Z-normalization**

1226 Z-normalization, also known as z-score normalization and
1227 “Normalization to Zero Mean and Unit of Energy” is a normalization method
1228 first mentioned by Goldin & Kanellakis (1995). The functionality of this method
1229 of time series normalization is to take the elements of an input vector and
1230 transform them into an output vector whose mean is approximately 0 with a
1231 standard deviation close to 1. The formula to achieve z-normalization is:

1232
$$\frac{x - \mu}{\sigma}$$

1233 Where x is an element within the time series, μ is the mean value within
1234 the time series and σ is the standard deviation of the time series.

1235 2. **Symbolic Aggregate ApproXimation (SAX)**

1236 This algorithm takes a z-normalized time series and transforms it into

symbolic representation to create words. Typically, any function using this method will have four parameters: The size of the window of time being observed, the length of the words that represent the data within that window, the size of the alphabet from which the words are built, and the type of numerosity reduction (P. Ordoñez et al., 2011). This data representation captures the shape of a time series, while also simplifying and facilitating pattern detection (J. Lin and Y. Li, 2009). For the purposes of implementation in this work, PAA size will be the equivalent of word size, size of the window of time being observed will be referred to as sliding window, and number of letters being used is alphabet size. In Figure 3, is a cartoon of how the transformation would look like for a sliding window of 180, PAA size of 7 and an alphabet size of 4. The time series in Figure 3 would be represented by the symbolic word CBDCCCC.

3. **Bag of words (BoW)**

Once the data is transformed, it is necessary to quantify it. BoW does this by quantizing each extracted word and then counting the frequency of each individual word contained in the time series (Senin & Malinchik, 2013). The final output of this is a table for each individual subject of study for whom we have a “count” of how many words represent them, and how often those words repeat. This by itself is not enough, as it does not highlight how much every word contributes to the overall shape of the time series.

4. **Term frequency–inverse document frequency (TF-IDF)**

This is a type of numerical statistic that measures how much a word

1260 contributes to the overall shape of the data (Neto. J, et al., 2000). Words that
1261 are unique within the data set are considered to weigh more and have a
1262 stronger effect on the shape of the data. For example, many bees in a hive
1263 may display highly rhythmic behaviours, but a group of them may only be
1264 active in the morning, another only in the afternoon, while a third group may
1265 be active all day long. The time series for a unique bee will be similar in shape
1266 to those who are active in the same shift. Therefore, all of the morning bees
1267 will have a distinct shape from the afternoon bees and in turn those two will
1268 have a different shape from the constant workers, even though all these bees
1269 may exhibit rhythmic behaviour.

1270 **Data Sets:**

1271 Once transformed using SAX, the data set was separated into the training and
1272 testing data sets. 90% of the individuals were randomly assigned to the training set,
1273 while the rest were used for the testing set.

1274 **Phase 4: Supervised Learning**

1275 Supervised machine learning is a type of machine learning where some of the
1276 data being analyzed already has the desired solutions (labels). The type of
1277 algorithms that use labels to learn an established pattern typically have one of two
1278 uses. One is to predict a target value based on a given set of numerical
1279 characteristics (features), this type of task is called regression. The other is the most
1280 common use for supervised learning, which is classification, or separating data into
1281 categories. In our case, the algorithm learns the patterns from the labeled data set

and uses this information to classify new data into the same categories (Géron, 2018).

5. **K Nearest Neighbors (KNN)**

The KNN algorithm, as used for classification purposes, is a non-parametric method that uses plurality to determine membership within a group (Bezdek et al., 1986). The user assigns a number K of minimum nearest neighbors that an individual being evaluated must have to be assigned a label (Bhatia, 2010) (Figure 4). The distances between neighbors are measured with Euclidean distances.

Decision Trees

Decision Trees consist of a unique central node that branches out into edges depending on the answer to binary questions. The binary questions represent a test on an attribute for a classification, each branch represents an outcome of the test which eventually ends in terminal nodes or leaves representative of the labels (Leonard, 2017). On the left-hand side on Figure 5 there is an example of how a decision tree could look like using our labels.

6. **Random Forest**

This algorithm is derived from Decision Trees. By definition, it is a combination of tree predictors, or the “forest”. Each tree depends on the values of a random vector sampled independently, and with the same distribution, for all other trees in the forest (Breiman, L. 2001). Based on a measurement of error, the best tree is chosen from the forest as the predictive model (Figure 5 right).

1305 All the code for phases 3 and 4 can be found in:

1306 <https://github.com/ComplejoC/CircadianSAX>.

1307 **Results:**

1308 **Clustering:**

1309 ***Clustering by Lomb Scargle Periodogram***

1310 An intuitive way to start is by separating individuals by the length of their day
1311 or periods. We reasoned that individual with similar lengths of day would cluster
1312 together. After doing consensus clustering, it was determined that 4 clusters was the
1313 optimal way to minimize errors and to have zero overlap between the clusters.
1314 Figure 6 is an illustration on how those clusters look like after using the PAM
1315 algorithm.

1316 Some clusters notably have more individuals than others, which illustrates
1317 that some mean lengths of days are more common than others. The requirements
1318 for cluster membership appear to be wider or shorter depending on the cluster. A
1319 closer inspection of the clusters, as illustrated in Figures 7 and 8, shows the
1320 characteristics of membership for each cluster.

1321 One may be tempted to look at all four graphs while looking for a pattern in
1322 clustering. Normally, the way to tell this story would be to only show the
1323 periodogram, because having all things together may seem confusing. Nevertheless,
1324 it's important to see all four graphs while validating the clusters for the biological
1325 interest, because that is how the circadian expert would evaluate the individuals.

1326 All these clusters in Figures 7 and 8 have in common the shape of the power
1327 spectrum, and not the length of their day per se. This leads to bees with different

1328 rhythms, but similarly shaped periodograms being put together. In some cases, like
1329 in Figure 7 B, the Lomb Scargle Periodogram could not determine a period for one of
1330 the individuals, but could determine it for the other. Furthermore, the period did not
1331 have to be similar for the individuals to group together. In some cases, we had
1332 individuals in the same cluster, for whom the difference in period would be 5 hours.

1333 ***Clustering by Autocorrelation***

1334 The next intuitive step is to cluster based on autocorrelation, as it is one of the
1335 most commonly used measures of rhythmicity. Because of its typical use, it stands to
1336 reason that the algorithm would cluster based on rhythm, which is what we are
1337 looking to describe. After using Consensus clustering, we determined that PAM with
1338 a K of 3 would work best to cluster the bees (Figure 9).

1339 We found no form of collision between the clusters (Figure 9). The dimensions
1340 of this plot do not appear to be representative of the data. Similarly, to a PCA, the
1341 sum of the percentages of both principal components should be as close to 100% as
1342 possible. There also seems to be an over-representation of individuals in cluster 1
1343 and few individuals both in cluster 2 and 3. Because it is known that this data set is
1344 diverse, it is worrying to see that most of the individuals grouped together.

1345 Contrary to our hypothesis, rhythmicity is not the feature on which the data is
1346 being clustered. It would appear that shape is once again the driving force behind
1347 our clusters. In the case of autocorrelation, it would appear to be the thickness of the
1348 autocorrelation plot.

1349 The second and third cluster appear to group things together by the thickness
1350 of the autocorrelation, where cluster 3 is thicker than 2. Nevertheless, contrary to

1351 what one would expect, it does not appear that all the individuals from these clusters
1352 have similar rhythms. To exemplify, in Figure 10. B., the thickness in the
1353 autocorrelation of both individuals is similar, but it's clear that the strength of their
1354 rhythms is contrasting.

1355 ***Clustering by Average Activity***

1356 For this group we choose to do clustering in a different manner. By using
1357 consensus, we determined that the K-means algorithm with a K of 3 was the best
1358 approach for this transformation (Figure 11. D.).

1359 The clusters for average daily activity are separated by magnitude of activity,
1360 with no consideration as to when the activity is being done. Neither of the locomotion
1361 plots (Figure 11. A. and B.) give the impression that any of the clusters have a
1362 common period of inactivity. The apparent lack of inactivity is the result of averaging
1363 individuals that are active in different times of the day. Figure 11. C illustrates that
1364 definitely, is the magnitude of activity that distinguishes all 3 clusters.

1365 **Classification:**

1366 ***K Nearest Neighbors (KNN):***

1367 To explore which transformation parameters optimized accurate classification,
1368 we explored all possible combinations of PAA size and Alphabet Size, as seen in
1369 Table 1, where PAA could be equal to: 3,4,6,8,12 or 24, and Alphabet Size could be:
1370 3,4,5,6 or 7. Additionally, we kept Sliding Window Size equal to 48, as we wanted
1371 our results to be analyzed in the circadian 24 hours. The best overall scores were
1372 achieved with a combination PAA of 4 and Alphabet Size of 3 (Table 1). Overall, no

1373 trend is evident when we experiment with the SAX parameters, i.e the alphabet size
1374 and PAA were typified by the combination of data set and classification algorithm.

1375 Often precision for the weakly rhythmic category was NA, which also reflected
1376 on the F1 measure. Contrastingly, the rhythmic category was the one which most
1377 often had any form of measurement, but it was the arrhythmic category that kept the
1378 better scores, although it has a larger amount of NAs than rhythm.

1379 ***Decision Trees and Random Forest:***

1380 With the intention of improving upon the results achieved from KNN, we
1381 tested the use of Decision Trees. After fitting the model using the rpart method from
1382 caret, we examined which Complexity Parameters were best for creating our
1383 predictions. Similarly, we build a Random Forest model using the same SAX
1384 parameters and the rf method for caret. In this case, we observed the number of
1385 randomly selected predictors to use in our prediction.

1386 Overall, both of the models are consistently better than KNN, and in one
1387 instance with PAA 6 and Alphabet 3, Random Forest performed the best with an
1388 accuracy of 0.889 (Table 2). Nevertheless, although both tree algorithms had better
1389 accuracy, whenever PAA or Alphabet size got larger than 4, the process became
1390 more computationally intensive. KNN on the other hand was not as computationally
1391 intensive, and on one occasion with PAA 4 and alphabet size of 3, achieved an
1392 accuracy of 0.8.

1393 Just like in KNN, tree methods had a large occurrence of NAs in precision,
1394 recall and F1 for the weakly rhythmic category. In Figure 12, we illustrate how those
1395 NAs happen. In the example, we are using the case of PAA = 4 and Alphabet Size =

1396 5 for KNN, but the same should apply for the other combinations of parameters and
1397 classification algorithms.

1398 We calculated accuracy by taking the sum of all correctly classified individuals
1399 from Figure 12.A. and divide them by the total number of individuals evaluated.
1400 Figure 12.C. illustrates how we calculate Precision, Recall and F-score. Precision
1401 was calculated by dividing all the True Positives (TP) by the sum of the TP and the
1402 False Positives (FP). In the example, for the WR category, both the total of TP and
1403 the sum of TP and FP equals zero, and therefore the calculation leads to an
1404 undefined value. Recall is the TP divided by the TP plus the False Negatives (FN).
1405 This is illustrated in our example for WR, where there is zero TP and four FN, and
1406 when plugged into our formula, it results in recall equal zero. Lastly, F-Score is
1407 calculated by doubling the product of Precision and Recall, and dividing it by the sum
1408 of Precision and Recall. This last calculation is highly dependent on the results of the
1409 two before it, and if both values equal zero, for example, or is even an undefined
1410 value, then F-score will not be defined. This is the case in those instances where the
1411 results for all possible evaluation metrics for one label is equal to NA.

1412 **Discussion:**

1413 **Clustering:**

1414 ***Clustering by Lomb Scargle Periodogram:***

1415 Using the Lomb Scargle periodogram to cluster our data resulted in clusters
1416 separated by the shape of the graph instead of the periodicity of each individual. This
1417 happened because periodograms are a power spectrum estimator and are usually
1418 used to evaluate the presence of oscillations within a dataset. While it does provide a

1419 measurement of periodicity, that is one isolated point within the periodogram, and we
1420 were providing to the clustering algorithms all of the data points within the graph.
1421 That is to say that in this case, more information about the shape and oscillation of
1422 the data did not, in fact, give a clear circadian set of instructions to the clustering
1423 algorithm. At least not in terms of the length of the biological clock.

1424 That being said, the clusters appear to have a loose sense of modality. In
1425 many cases like in Figure 7.A., the number of curves in the periodogram seem to be
1426 reflective of the number of peaks of activity of the bee. This sense of shape is likely
1427 what is informing the clusters and should be explored further with proper parameter
1428 tuning for the Lomb Scargle transformation.

1429 ***Clustering by Autocorrelation:***

1430 Clustering by autocorrelation resulted in clusters that were either informed by
1431 the thickness of the autocorrelation or by simply not having a discernible pattern. The
1432 autocorrelation has its usefulness to circadian science in calculating period and
1433 rhythmicity. Knowing this, one would expect that rhythm would be the main informing
1434 feature to the clustering algorithm in this experiment. However, in reality, what best
1435 informed the composition of the clusters was once again the shape of the data. This
1436 is why there exist individuals within groups with different rhythmicities, but with a
1437 similar shape to their activity pattern.

1438 ***Clustering by Average Activity:***

1439 In this last attempt, we not only managed to cluster in such a way that no
1440 collision was detected, but additionally, the characteristics of each cluster was clear
1441 and easy to define. Unfortunately, even though the individuals clustered

1442 appropriately by average activity, the groups did not share circadian similarities,
1443 which was our goal. For example, the individuals in any given cluster may have the
1444 same amount of activity over time, but did not necessarily share the same length of
1445 day (period) or even be active during the same hours of the day, among other key
1446 aspects.

1447 For all three univariate clustering, successful clustering does not equate to
1448 significant positive results in terms of a practical question. The fact that all three
1449 clustering applications did not yield a circadian significance does not imply a failure
1450 in part of the algorithms, but instead, it shines a light on how complex time series
1451 questions can be, and more so if we take the biological significance into
1452 consideration. The way a circadian scientist would divide individuals in a population
1453 into discrete groups is a multivariate process that would take into consideration all of
1454 the transformations of the data that we evaluated individually with clustering.
1455 Although other types of transformations could be considered for univariate clustering,
1456 a multivariate approach should also be considered. However, multivariate clustering
1457 is not a trivial pursuit. Using transformations that can inform the shape and intensity
1458 of the activity may definitely inform a good multivariate approach but developing the
1459 mathematical and architectural tools necessary for this is a thesis in itself. This is
1460 due to the size of time series data. We transform the data to reduce dimensionality,
1461 but that does not mean we eliminate the continuous time component when we
1462 transform, in fact we accentuate it. To reduce dimensionality further for the sake of
1463 clustering may do an injustice to the data set. Therefore, the best way to group this
1464 type of data set is one that conserves the fidelity of the data while also minimizing

1465 the amount of computational resources necessary to achieve the task, which is why
1466 we did classification.

1467 **Classification:**

1468 When we validated the results of the clustering analysis, we noticed a number
1469 of patterns that appeared to be of circadian nature. Because multivariate clustering is
1470 a more complex problem than what we were equipped to handle, these patterns
1471 offered an alternative for analyzing these data. We created a user defined set of
1472 labels and had a group of experienced and naive users assigned those labels to the
1473 data set.

1474 For all three classification algorithms, the weakly rhythmic category was the
1475 one that caused the most difficulties. Consistently, it was the category that most
1476 often returned NA in the evaluation metrics. This suggests that we should re-
1477 evaluate what makes an individual weakly rhythmic, or even subdivide it into smaller
1478 categories still. On the other hand, we could also consider that none of the three
1479 algorithms used may be appropriate for the type of data we are using. Nevertheless,
1480 for the Rhythmic and Arrhythmic categories, all three classification methods
1481 performed adequately, which leads us to believe that more appropriate labels are
1482 needed.

1483 ***K Nearest Neighbors (KNN):***

1484 Our results suggest that combinations of smaller PAA size and Alphabet size
1485 parameters in SAX transformation yield better results when evaluating the model.
1486 More often than not, KNN did not return a model with accuracy higher than 0.5,

1487 which suggests that the use of this classification algorithm requires fine-tuning for it
1488 to return correctly classified results.

1489 ***Decision Trees and Random Forest:***

1490 In contrast, both tree algorithms consistently got an accuracy over 0.5,
1491 although only once did they outperform the highest KNN result. The optimal
1492 parameters for the best model using Decision Trees follows the tendency in the
1493 literature of smaller PAA and Alphabet size being the most optimal (J. Lin, 2007).
1494 Although one could argue that PAA 6 and Alphabet 3 are still relatively small. The
1495 better model came at the cost of computational power and took considerably longer
1496 to build.

1497 **Conclusions:**

1498 Clustering, although a good first step to gaining intuition for the behavior of
1499 the data, does not result in strong conclusions. For all three attempts, the algorithm
1500 used the parameter given in a way we did not expect. For example, for periods using
1501 the shape of the data rather than its values. While these experiments gave us a good
1502 intuition for the data set, the ultimately did not satisfy our need of grouping bees by
1503 circadian phenotypes. The next steps would be to attempt clustering with other
1504 transformations or even a multivariate analysis. Because our goal was to facilitate
1505 circadian analysis, we decided it would be simpler to group using methods that
1506 mimic how a circadian expert evaluates multiple parameters of an individual to
1507 describe its characteristics.

1508 In our classification experiments for all possible combinations of parameters
1509 evaluated, the Weakly Rhythmic label caused the most difficulty in classification.

1510 Often returning zeros and NAs for precision and recall. At this moment, we have not
1511 noticed any common patterns in the WR classifications. Nevertheless, this persistent
1512 difficulty in classification does inspire a reconsideration on how we are labeling these
1513 individuals. Because the individuals are not consistently being misclassified as
1514 rhythmic or arrhythmic, it stands to reason that the errors are not necessarily caused
1515 by the users mislabeling the individuals, but that just using three generalized labels
1516 could be obfuscating certain behaviors in a category. Therefore, experimenting with
1517 separating the weakly rhythmic category into further categories may facilitate the
1518 classification process. Furthermore, these more in-depth classifications can help
1519 reflect the biological reality of *L. malachurum*.

1520 All is not lost, as we did build two different classification models with 0.80
1521 accuracy or more. This demonstrated that at least two of the categories built with
1522 user input were adequate to inform a model. Of the two models The KNN one paired
1523 with SAX transformation parameters of PAA 3 and alphabet size 4 (Table 1) is the
1524 less accurate at 0.80, nevertheless it was the one that ran the fastest and used the
1525 least amount of computation resources. While the Decision Tree model with PAA 6
1526 and Alphabet size 3 (Table 2) had an accuracy at 0.89 but took longer to run and
1527 more computational resources. Therefore, the choice of the best model is arguably
1528 determined by the resources available to the user and having more than one model
1529 from which to choose is beneficial for those who may not have much computing
1530 power available to them.

1531 Ultimately, we successfully set the basis for an evaluation pipeline for
1532 circadian data that isn't heterogeneous. This will undoubtedly facilitate the evaluation

1533 of organisms that naturally express multiple phenotypes of circadian rhythms. The
1534 next steps are to test this pipeline with other organisms to observe whether or not it
1535 is indeed generalizable.

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561 **Figure Legends:**

1562 **Figure 1: Locomotor activity over time data from *L. malachurum* displays**

1563 **heterogeneity.** All the 3 figures contain a Double plotted Actogram, a Daily Average
1564 activity plot, a Lomb Scargle Periodogram and an Autocorrelation. A double plotted
1565 actogram is a visualization typically used to represent the Circadian locomotor
1566 activity rhythms. The daily average activity plot is representative of the average
1567 activity done by each individual for the duration of the study. Lomb Scargle
1568 periodogram is an algorithm that can be used to describe the period of unevenly
1569 sampled time series data sets. The autocorrelation of a time-series measures how
1570 similar a time-series is with a forward or backwards shifted version of itself. In **A**, **B**,
1571 and **C**, is a sample of the diversity of behaviours displayed in the data set.

1572 **Figure 2: Average of all 98 individuals in the dataset is not representative of**
1573 **any one individual due to heterogeneity.** The single plotted actogram illustrated
1574 was created by taking the average of the activity for all 98 individuals done in the
1575 four day observational period.

1576 **Figure 3 A mock-up of a normalized time series transformed with SAX.** A sliding
1577 window of 180, PAA size of 7 and an alphabet size of 4 together make the symbolic
1578 word CBDCCCC.

1579 **Figure 4: A mock-up of how KNN works.** The pentagon is a new piece of data
1580 plotted into n-dimensional space here simplified as two-dimensional space. By
1581 consensus, the pentagon is closest to the group labeled AR and therefore will be
1582 classified as such.

1583 **Figure 5 Mock-up illustrating Decision Trees (left) and Random Forest (right).**

1584 For Decision trees only one chart is considered as a classifier, whereas in Random

1585 Forest many trees are considered and the best tree is chosen based on a

1586 measurement of error

1587 **Figure 6: PAM with $K = 4$ for the Lomb Scargle Periodogram resulted in highly**

1588 **representative discrete clusters.** We used principal components to plot the

1589 clusters, where the X and Y Axes are the principal components of the data set

1590 **Figure 7: Clusters resulting from L. S are separated by the shape of the**

1591 **periodogram Part 1.** In green A) are examples of membership from cluster 1,

1592 characterized by two peaks in the periodogram that cross the horizontal line or the

1593 first barely does. In orange B) are examples of membership from cluster 2, where

1594 the individuals have no peak, or if they have one, it barely touches the horizontal line

1595 of the periodogram

1596 **Figure 8: Clusters resulting from L. S are separated by the shape of the**

1597 **periodogram Part 2.** In purple A) are examples of the membership in cluster 3.

1598 Where individuals either have two peaks in the periodogram, where the first is far

1599 from touching the horizontal line, or just the one peak. In fuchsia B) are examples of

1600 membership of cluster 4. The individuals in this group have two peaks in their

1601 periodograms, but the first peak is small and sometimes unstable in terms of shape.

1602 **Figure 9: PAM with $K = 3$ for Autocorrelation coefficient clusters of low**

1603 **representation power.** Utilizing the same form of dimensionality reduction in Figure

1604 3, we plotted the clusters resulting from PAM. Most of the individuals in the data set

1605 are clustered into cluster 1, while the rest divided into the other two clusters.

Figure 10: PAM with K =3 for Autocorrelation coefficient clusters by thickness

of Autocorrelation. In green, A) are examples of the membership in cluster 1.

Individuals in this cluster have a thinner autocorrelation graph than that of clusters 2 and 3, passing the horizontal line at most once. Nevertheless, the thickness does not seem to be consistent across the membership of cluster 1. In orange, B) are examples of the membership in cluster 2, where the individuals observed pass the horizontal line more than once and show larger density than cluster 1, but less than cluster 3. In purple, c) are examples of the membership in cluster 3, where the individuals observed have the thickest autocorrelation plot.

Figure 11: K-means with K= 3 for Average Daily Activity clusters by frequency.

A) and B) show the shape of the activity contained within the clusters. In A), the graph shows the clusters individually, while B) shows them together. C) Illustrates the distribution of average activity within each cluster. D) Principal components illustration of the clusters.

Figure 12: Prevalence in NAs is due to poor consistent classification. Here we

illustrate the process of calculating Accuracy, Precision, Recall and F-score, all measurements that we used to test how good our model is. These are the values taken from Table 1 for PAA = 4 and alphabet size = 5. On A is the test labels matched to the predicted labels, in red are the incorrectly classified individuals and in blue are the correctly classified individuals. In B. are the calculations for accuracy, which is the sum of all correctly classified individuals divided by the total of individuals, in addition there is also a confusion matrix. True positives (TP) are the values within each cell or the correctly classified individuals, false positives (FP) and

false negatives (FN) are the incorrectly classified individuals, which are indirectly viewed in the sum of columns or rows. Lastly, in C., a table illustrating the calculations for the rest of the measurements were Pre = Precision and Re = Recall.

Tables and Figures:

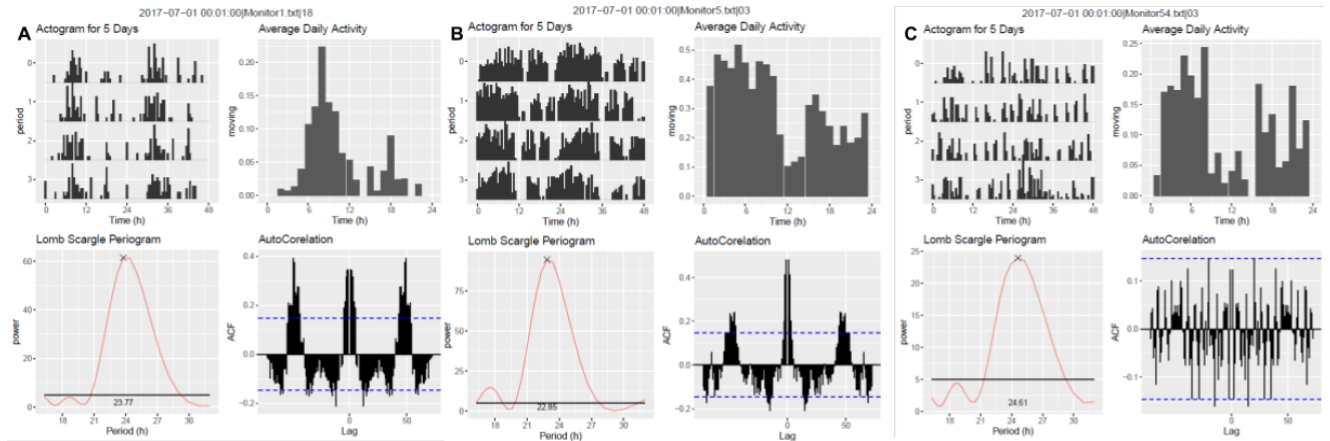
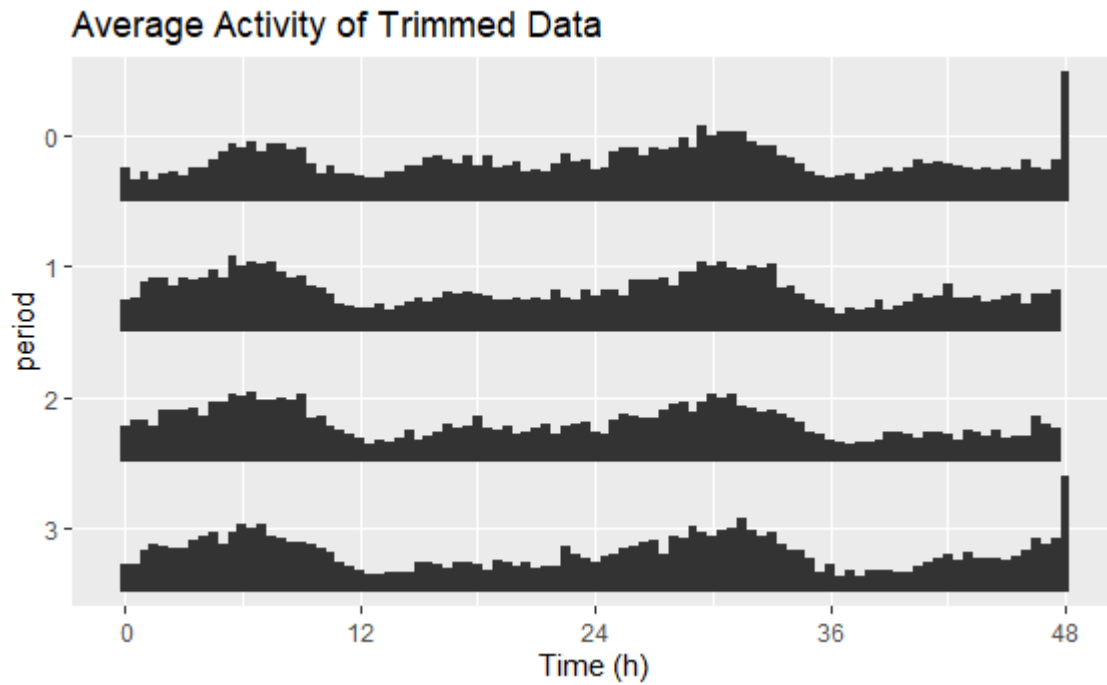


Figure 1: Locomotor activity over time data from *L. malachurum* displays

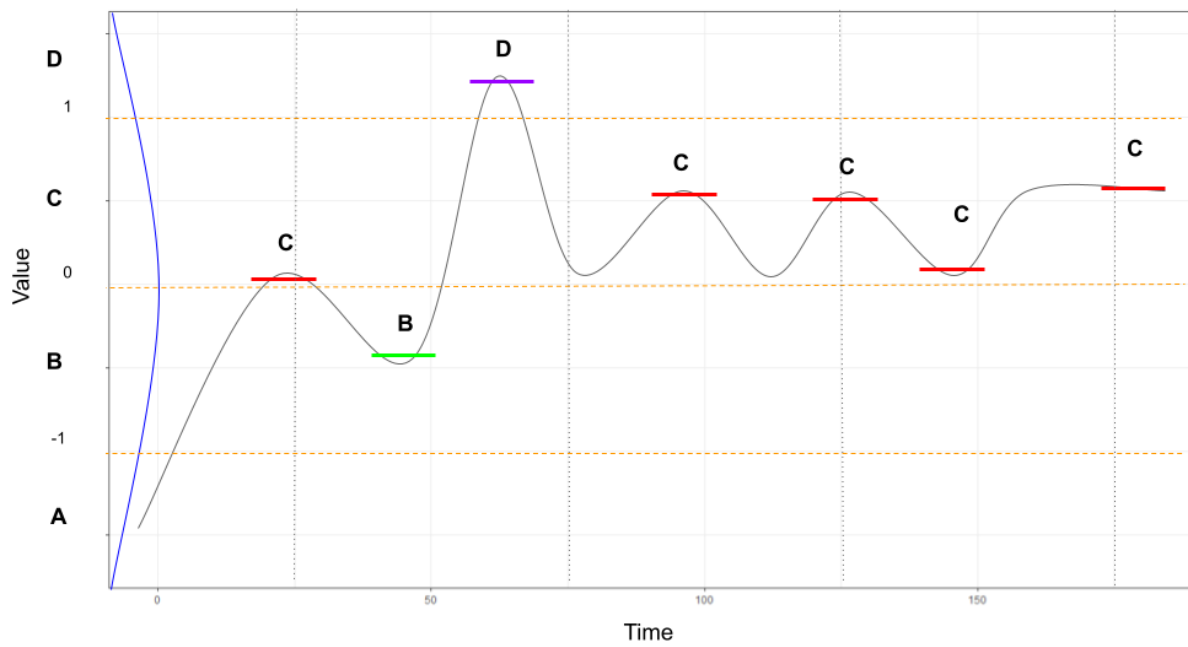
heterogeneity.



1636

1637 **Figure 2:** Average of all 98 individuals in the dataset is not representative of any one
1638 individual due to heterogeneity.

1639



1640

1641 **Figure 3** A mock-up of a normalized time series transformed with SAX.

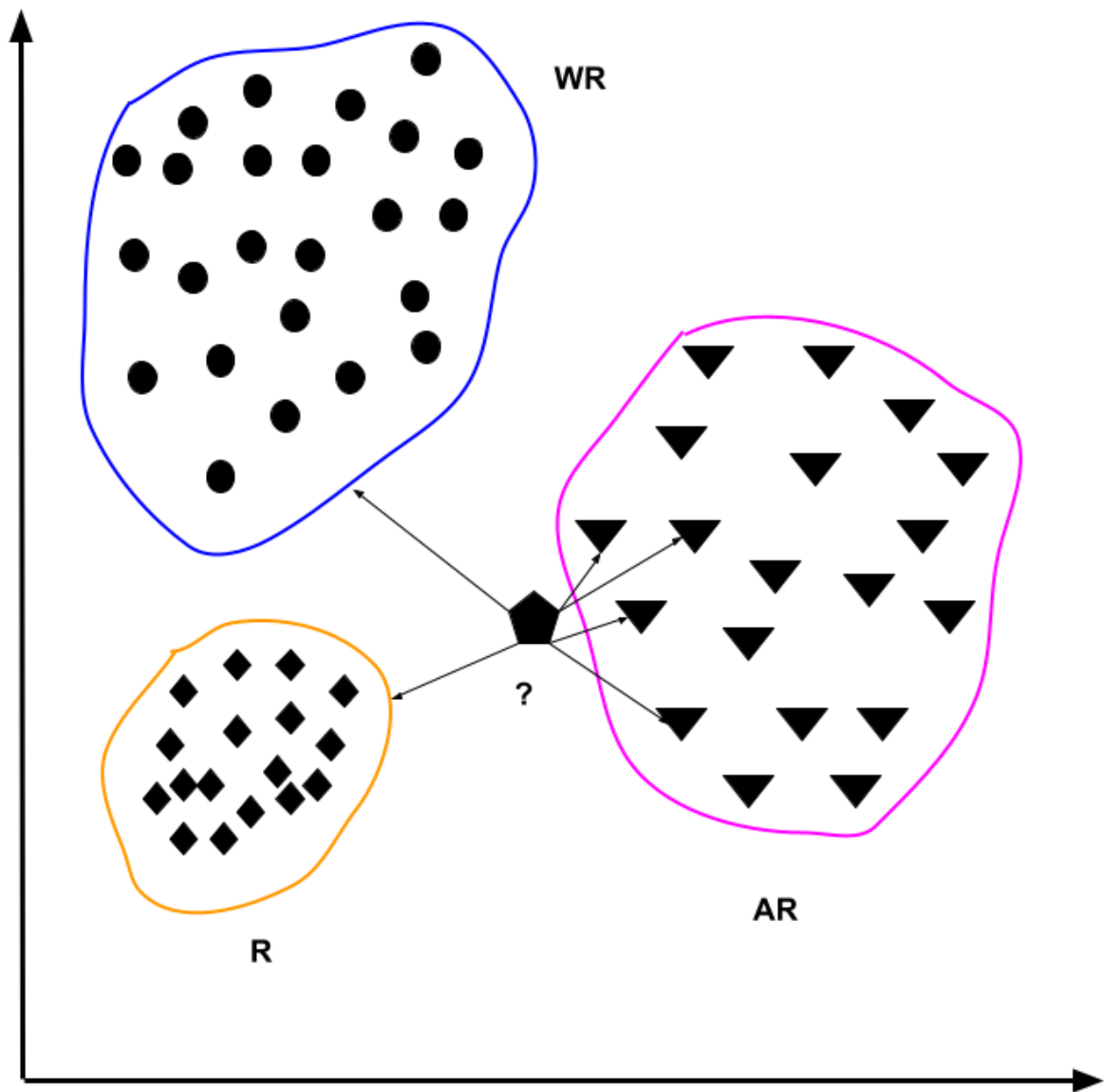


Figure 4: A mock-up of how KNN works.

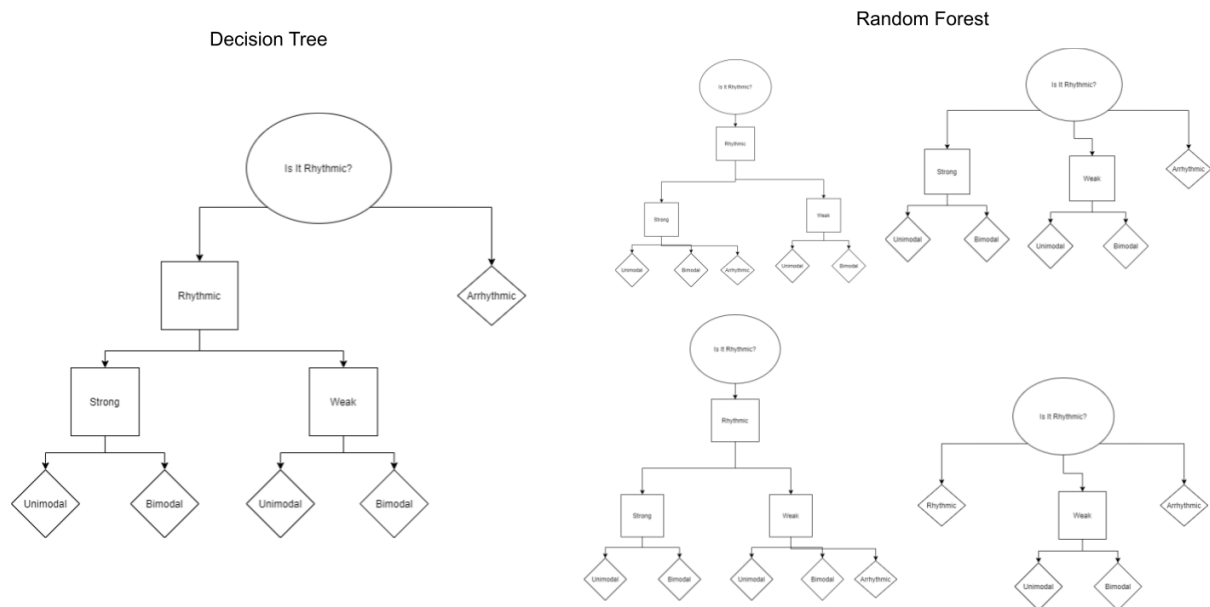


Figure 5 Mock-up illustrating Decision Trees (left) and Random Forest (right).

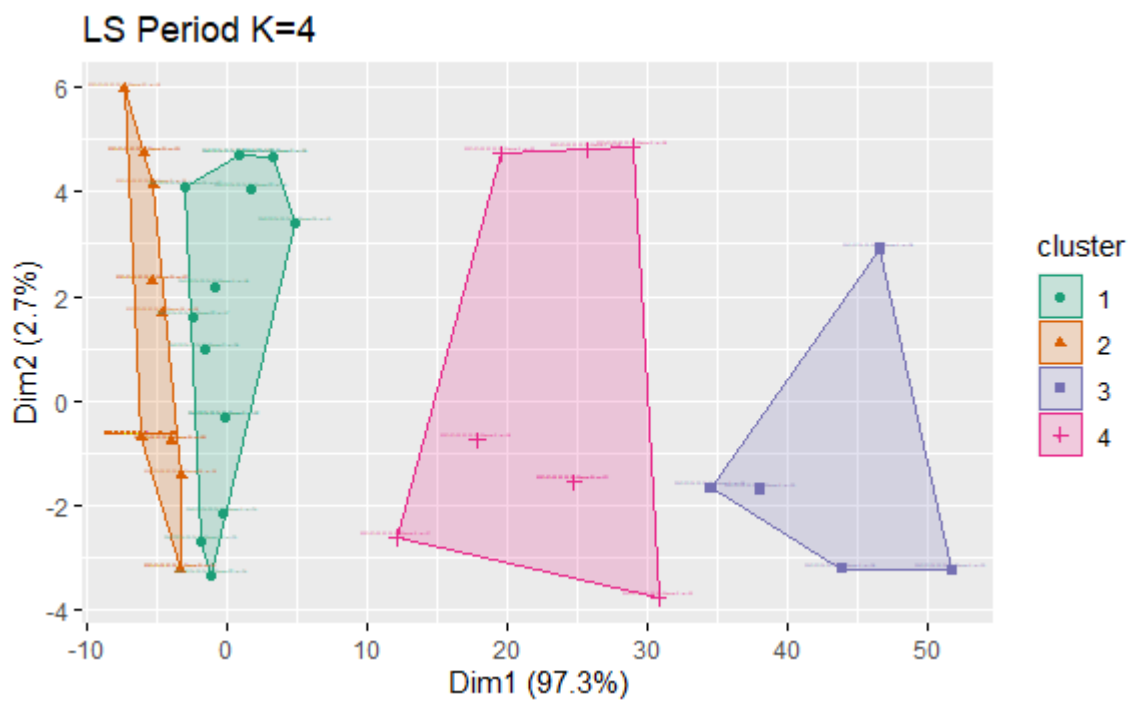


Figure 6: PAM with $K = 4$ for the Lomb Scargle Periodogram resulted in highly representative discrete clusters.

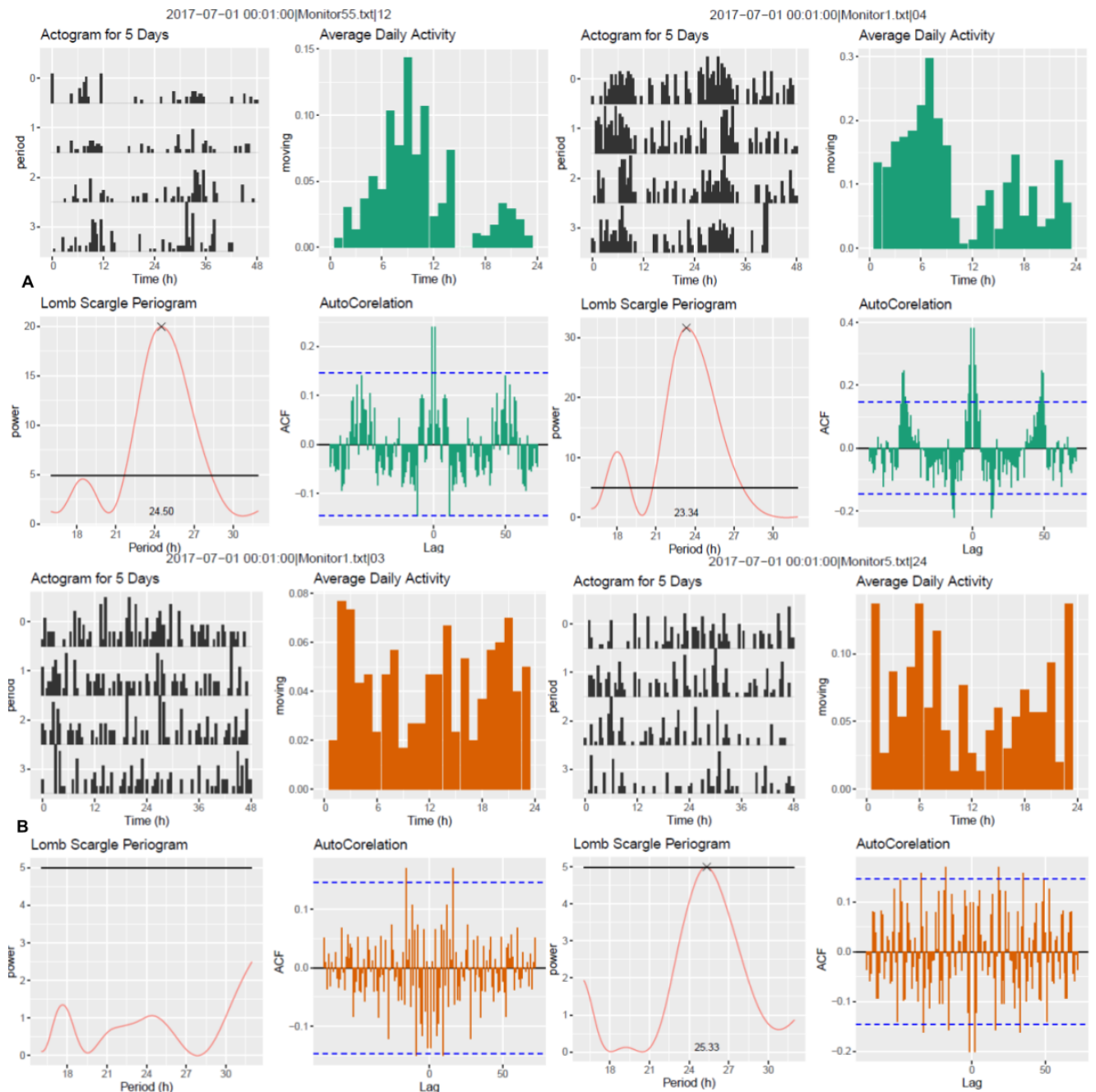


Figure 7: Clusters resulting from L. S are separated by the shape of the periodogram Part 1.

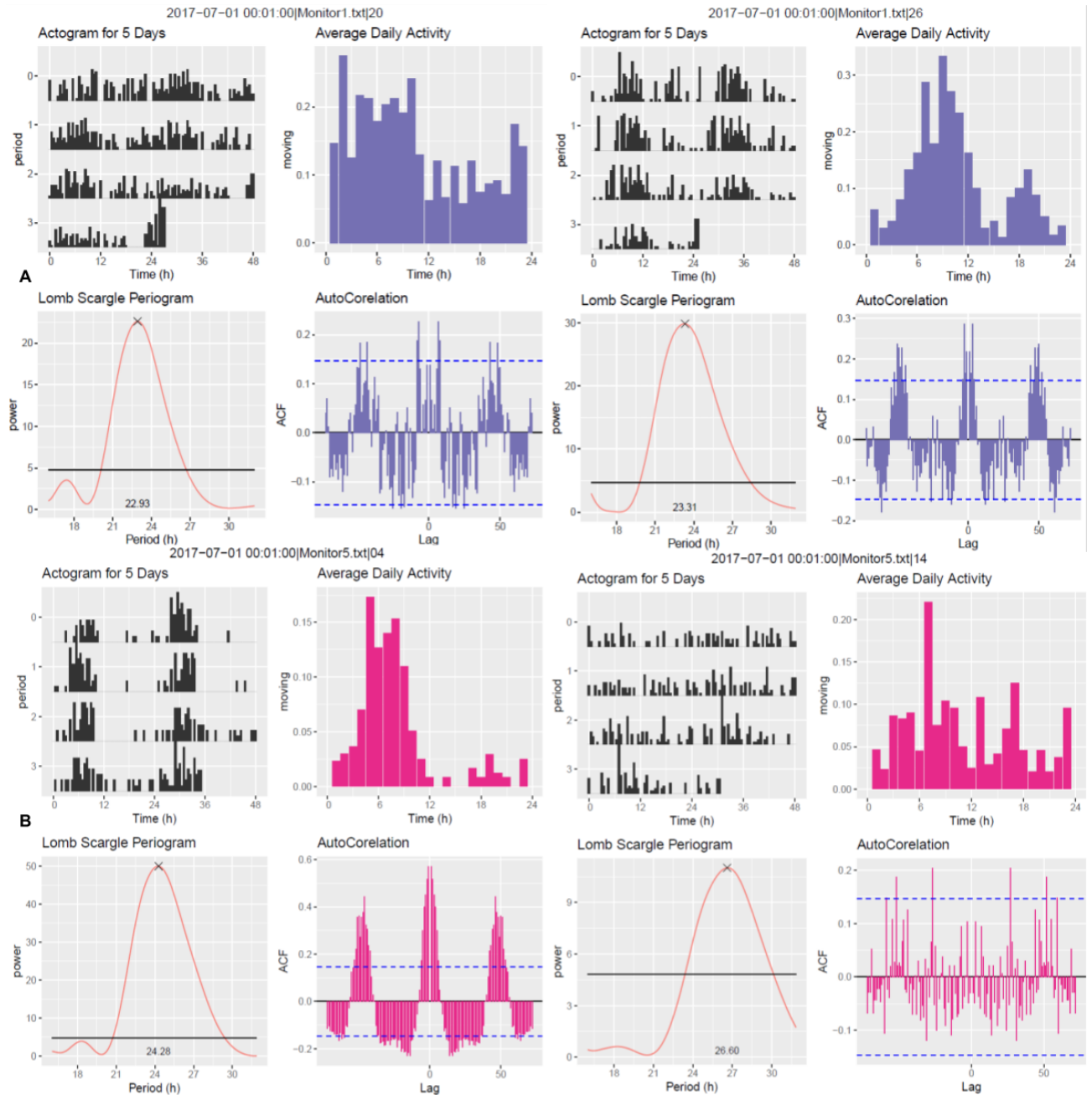
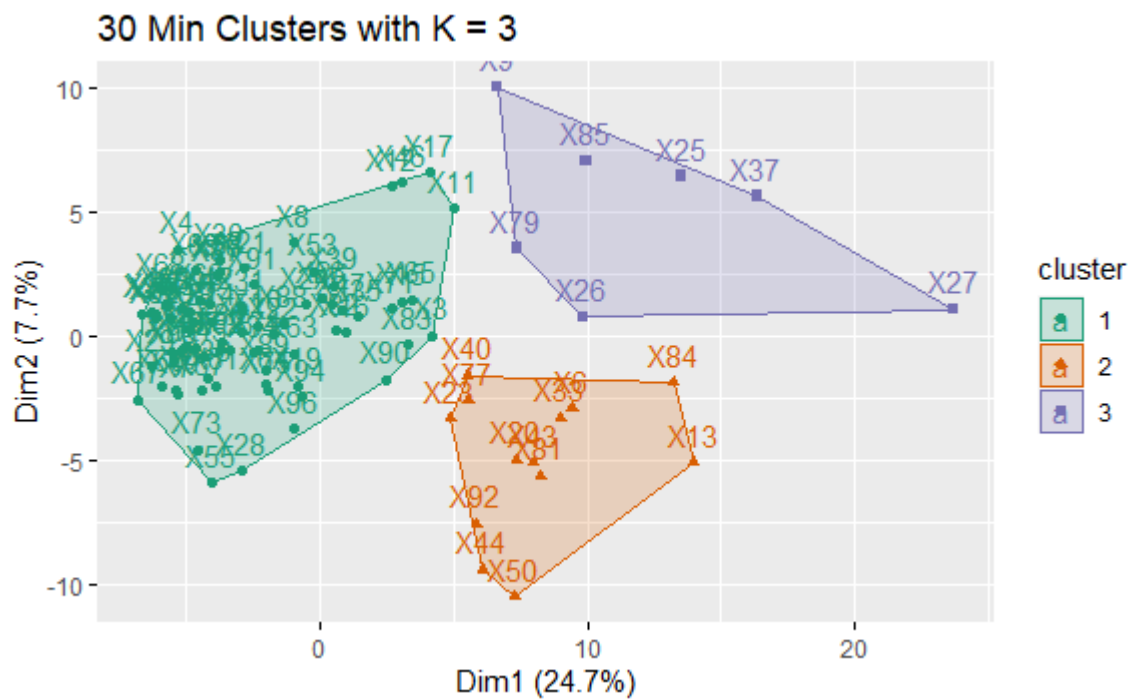


Figure 8: Clusters resulting from L. S are separated by the shape of the periodogram Part 2

1656



1657

1658 **Figure 9:** PAM with K =3 for Autocorrelation coefficient clusters of low representation
1659 power.

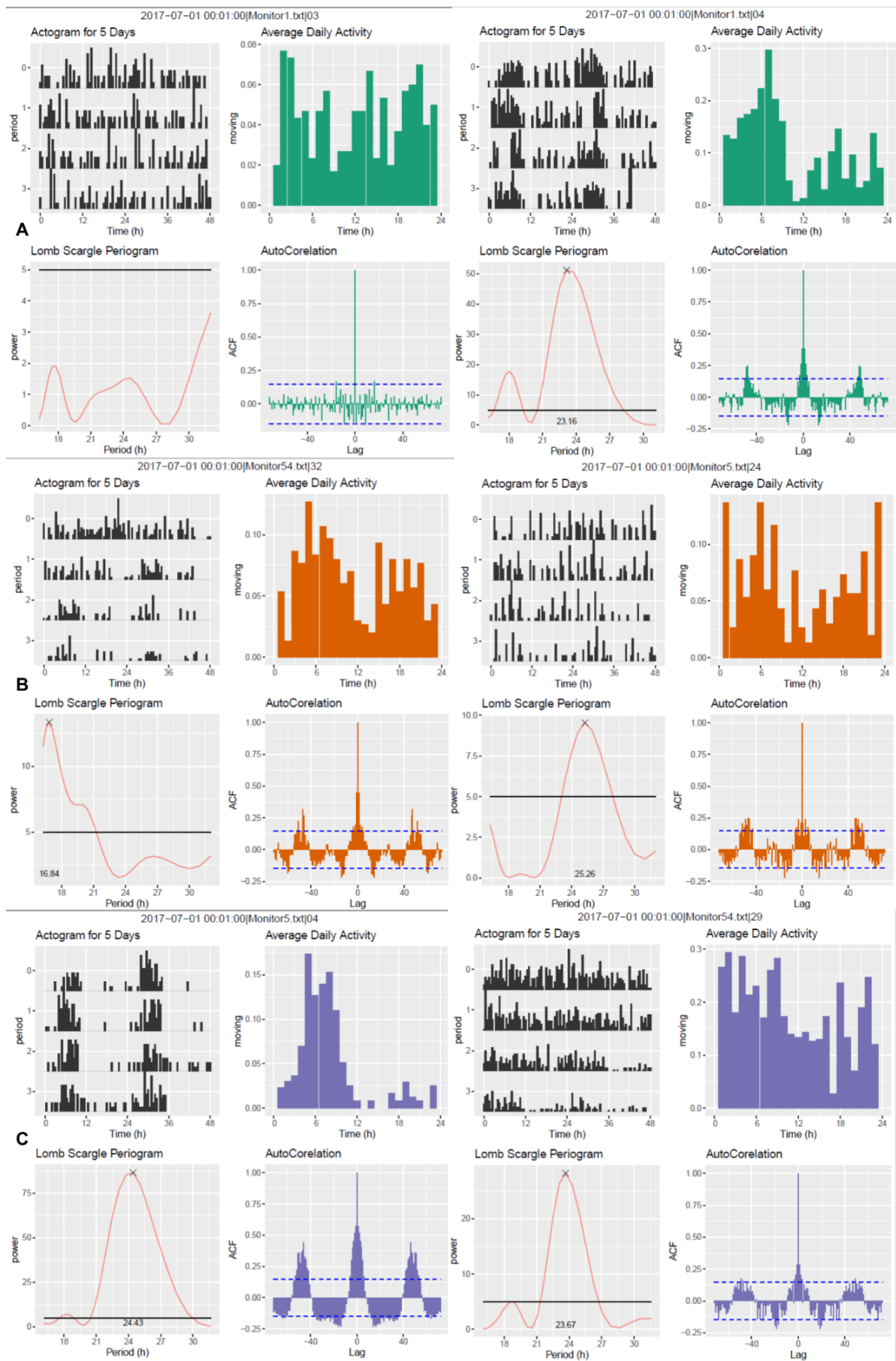


Figure 10: PAM with K =3 for Autocorrelation coefficient clusters by thickness of Autocorrelation.

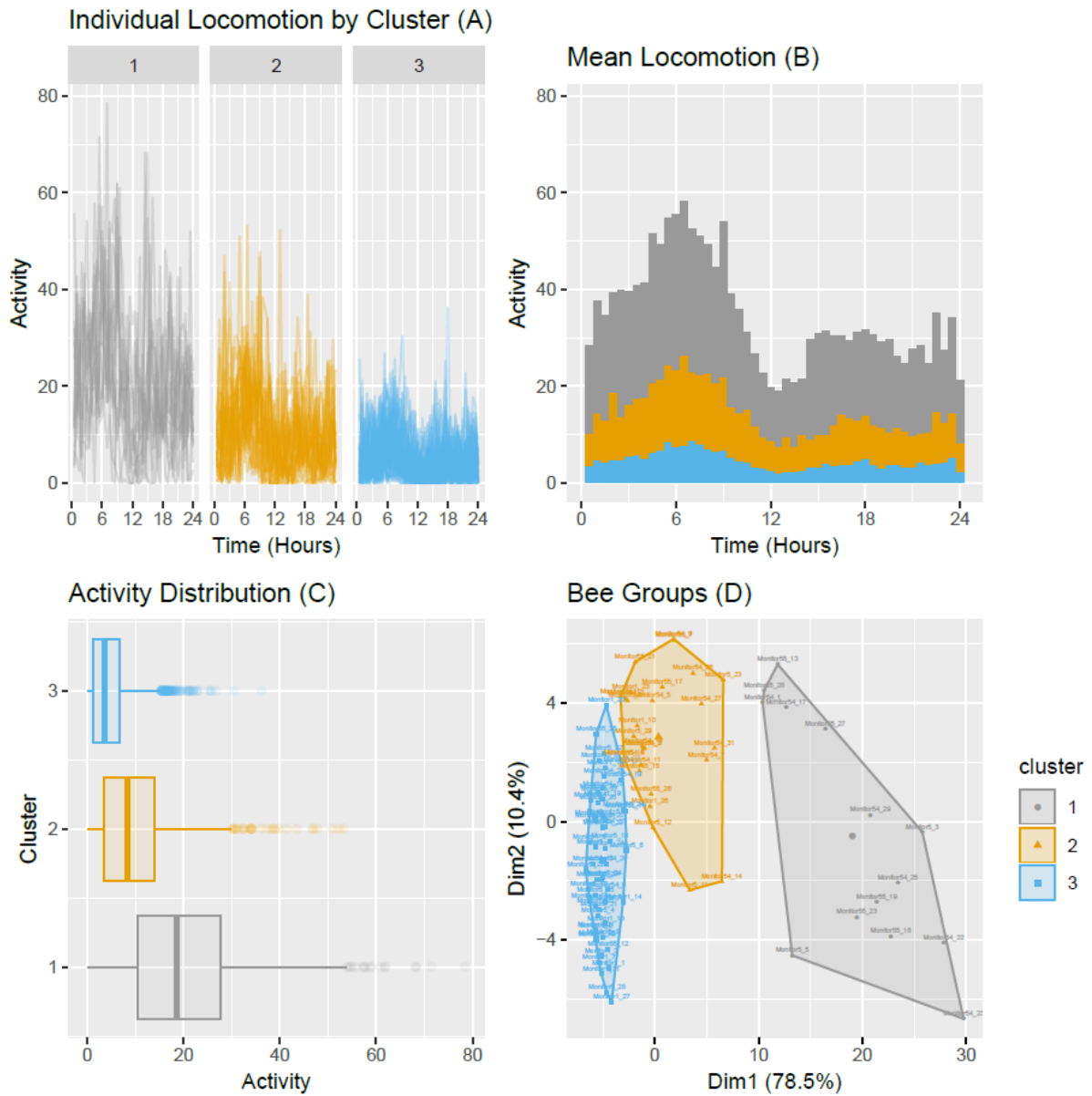


Figure 11: K-means with K= 3 for Average Daily Activity clusters by frequency.

1668 **Table 1:**Results of KNN by varying PAA Size (PAA) and Alphabet Size (Alphabet).

1669 Sliding Window (SW) was continually equal to 48.

Parameters			Precision			Recall			F1			Accuracy
PAA	Alphabet	SW	R	WR	AR	R	WR	AR	R	WR	AR	
3	3	48	0.5	0	0.25	1	0	0.25	0.67	NA	0.25	0.30
3	4	48	1	0	0.375	0.5	0	0.75	0.5	0.67	NA	0.40
3	5	48	0.5	0	0.4	1	0	0.5	0.67	NA	0.4	0.40
3	6	48	0.5	NA	0.5	1	0	0.75	0.67	NA	0.6	0.50
3	7	48	1	0.67	0.5	0.5	0.5	0.75	0.67	0.57	0.6	0.60
4	3	48	0.67	0.75	1	1	0.75	0.75	0.8	0.75	0.86	0.80
4	4	48	1	0.5	0.6	1	0.25	0.75	0.67	0.4	0.67	0.60
4	5	48	0.67	NA	0.57	1	0	1	0.8	NA	0.73	0.60
4	6	48	1	1	0.57	0.5	0.5	1	0.67	0.67	0.2	0.70
4	7	48	0.3	0	0.5	0.5	0	0.75	0.4	NA	0.6	0.40
6	3	48	0.5	NA	0.5	0.5	0	1	0.5	NA	0.67	0.50
6	4	48	0.33	0.33	0.5	0.5	0.25	0.5	0.4	0.29	0.5	0.40
6	5	48	1	NA	0.5	1	0	1	1	NA	0.67	0.60
6	6	48	NA	NA	0.4	0	0	1	NA	NA	0.57	0.40
6	7	48	0.67	NA	0.57	1	0	1	0.8	0.73	NA	0.60
8	3	48	NA	0.4	NA	0	1	0	NA	0.57	NA	0.40
8	4	48	0.26	NA	0.67	1	0	0.5	0.4	NA	0.57	0.40
8	5	48	0.17	NA	0.75	0.5	0	0.75	0.25	NA	0.75	0.40
8	6	48	NA	NA	0.4	0	0	1	NA	NA	0.57	0.40
8	7	48	0.5	0.5	0.57	0.5	0.25	1	0.5	0.4	0.73	0.60
12	3	48	NA	NA	0.4	0	0	1	NA	NA	0.57	0.40
12	4	48	0.2	NA	NA	1	0	0	0.33	NA	NA	0.20
12	5	48	0.2	NA	NA	1	0	0	0.33	NA	NA	0.20
12	6	48	0.25	0.5	0.5	0.5	0.5	0.25	0.33	0.5	0.33	0.40
12	7	48	0.2	NA	NA	1	0	0	0.33	NA	NA	0.20
24	3	48	0.25	0.5	0.5	0.5	0.5	0.25	0.33	0.5	0.33	0.40

24	4	48	0.2	NA	NA	1	0	0	0.33	NA	NA	0.20
24	5	48	0.2	NA	NA	1	0	0	0.33	NA	NA	0.20
24	6	48	0.2	NA	NA	1	0	0	0.33	NA	NA	0.20
24	7	48	0.2	NA	NA	1	0	0	0.33	NA	NA	0.20

1670

1671 **Table 2:** Results Trees by varying PAA Size (PAA) and Alphabet size (Alphabet).

1672 Sliding Window (SW) was continually equal to 48.

Parameters			Precision			Recall			F1			Accuracy	Algorithm
PAA	Alphabet	SW	R	WR	AR	R	WR	AR	R	WR	AR		
3	3	48	1	NA	0.5	0.75	0	1	0.8571	NA	0.667	0.667	DT
3	3	48	1	0.5	0.667	1	0.5	0.667	1	0.5	0.667	0.778	RF
4	3	48	1	NA	0.5	0.75	0	1	0.75	NA	0.667	0.667	DT
4	3	48	1	0	0.5	1	0	0.667	1	NA	0.571	0.667	RF
6	3	48	1	NA	0.6	1	0	1	1	NA	0.75	0.778	DT
6	3	48	1	1	0.75	1	0.5	1	1	0.667	0.8571	0.889	RF
3	4	48	1	NA	0.4286	0.5	0	1	0.667	0.6	NA	0.556	DT
3	4	48	1	NA	0.5	0.75	0	1	0.8571	0.667	NA	0.667	RF
4	4	48	0.8	NA	0.5	1	0	0.667	0.889	NA	0.5714	0.667	DT
4	4	48	0.8	NA	0.5	1	0	0.667	0.889	NA	0.5714	0.667	RF
6	4	48	0.5	NA	0	1	0	0	0.667	NA	NA	0.444	DT
6	4	48	1	NA	0.6	1	0	1	1	NA	0.75	0.778	RF
3	5	48	1	NA	0.4286	0.5	0	1	0.667	NA	0.6	0.556	DT
3	5	48	1	1	0.6	0.75	0.5	1	0.8571	0.667	0.5	0.778	RF
4	5	48	1	1	0.6	0.75	0.5	1	0.8571	0.667	0.75	0.778	DT
4	5	48	1	0.5	1	0.75	1	0.667	0.8571	0.667	0.8	0.778	RF
6	5	48	0.6	NA	0.75	0.75	0	1	0.667	NA	0.8571	0.667	DT
3	6	48	0.667	NA	0.667	1	0	0.667	0.8	NA	0.667	0.667	DT
3	6	48	1	0	0.5	0.75	.	0.667	0.8571	NA	0.5714	0.556	RF

4	6	48	0.8	NA	0.5	1	0	0.667	0.889	NA	0.5714	0.667	DT
4	6	48	0.8	NA	0.5	1	0	0.667	0.889	NA	0.5714	0.778	RF
3	7	48	0.8	NA	0.75	1	0	1	0.889	NA	0.8571	0.778	DT
3	7	48	1	NA	0.6	1	0	1	1	NA	0.75	0.778	RF
4	7	48	0.8	NA	0.75	1	0	1	0.889	NA	0.8571	0.778	DT
4	7	48	1	NA	0.6	1	0	1	1	NA	0.75	0.778	RF
3	8	48	1	NA	0.6	1	0	1	1	NA	0.75	0.778	DT
3	8	48	1	NA	0.6	1	0	1	1	NA	0.75	0.778	RF
4	8	48	0.75	NA	0.4	0.75	0	0.667	0.75	NA	0.5	0.556	DT
4	8	48	1	Na	0.6	1	0	1	1	NA	0.75	0.778	RF

1673

PAA = 4, Alphabet size = 5

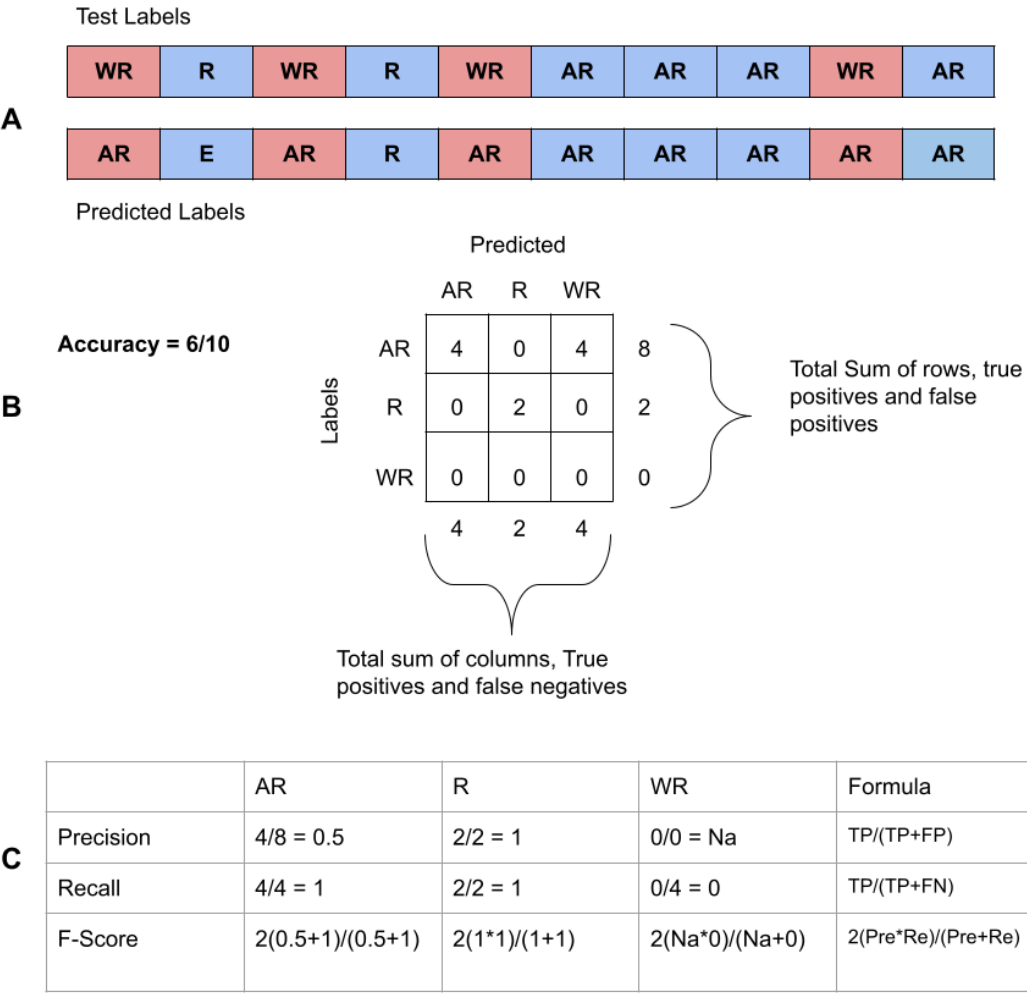


Figure 12: Prevalence in NAs is due to poor consistent classification.

1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707

Discussion/Conclusion

1708 Discussion/Conclusion

1709 When describing the phenotypic expression of circadian rhythms in our four
1710 Halictid bees, we noted that rhythm variability is parallel to that of levels of sociality.
1711 From least to most varied in behaviors: *S. curvicornis*; Solitary and with one activity
1712 pattern, *L. ferreirii*; Communal with two activity patterns, *L. enatum*: Primitively
1713 eusocial with three activity patterns and *L. malachurum*; Facultatively eusocial with
1714 five activity patterns. Out of the four species evaluated, *S. curvicornis* is the only
1715 specialist, and its relationship to *Campanula arvensis* could explain its rigorous
1716 biological clock. However, within the context of this work, that explanation falls short,
1717 because all three of the *Lasioglossum* species are generalists, yet they express a
1718 significant difference across their circadian parameters, not to mention that *L.*
1719 *malachurum* was caught in the same time place and flower as *S. curvicornis*, yet
1720 their patterns for daily locomotion are polar opposites. This is similarly reflected in *L.*
1721 *enatum* and *L. ferreirii* sharing the same exact niche and still displaying a different
1722 number of activity patterns. It stands to be seen if these results hold up when
1723 considering sexual dimorphism and seasonality. Notwithstanding, our observations
1724 set the foundation for asking more complex questions about the influences of
1725 sociality in the expression of circadian rhythms in Hymenopterans. It is our belief that
1726 these observed differences in rhythm are related to shift work and that understanding
1727 circadian rhythms may give a strong insight into the mechanisms that lead to the
1728 evolution of complex social organizations. Therefore, we find it worthwhile to look
1729 further into this relationship and identify if there is a correlation between the task
1730 being done by a bee and their circadian phenotype.

1731 While describing the *Lasioglossum* species, we created a labeling scheme
1732 that highly resembles a classification decision tree. Similarly, the original process for
1733 exploring the reproducibility of the categories observed in the data set was
1734 reminiscent of KNN, where a consensus of users separated the data into discrete
1735 pre-established groups. We divided the data this way, because the time series
1736 behavior being displayed by the *Lasioglossum* bees was highly heterogeneous. To
1737 streamline the *a priori* grouping process, we used clustering analysis as well as
1738 classification on the *L. malachurum* dataset. A systematic and replicable
1739 methodology for preprocessing locomotor activity data will not only make data
1740 analysis faster and easier, but it will also strengthen the reliability of the results.

1741 Clustering proved to be a successful method for grouping individual bees.
1742 However, the groups were not of circadian significance, as they were not grouped by
1743 rhythm. Nevertheless, the process was pivotal in creating a deeper understanding of
1744 the behavior of the data set. Classification, on the other hand, was a complete
1745 success. We transformed the data using SAX with the intent of reducing the
1746 dimensions of the data set while still keeping its shape. Once the data was
1747 transformed, we applied three different classification algorithms, with decision trees
1748 being the algorithm that best classified the data, KNN achieving the second-best
1749 results, and random forest coming in last place. These findings are reflective of how
1750 intuition gained by understanding a data set is the most valuable step in any ML
1751 workflow. Using methods that approximate how a user may separate data into
1752 discrete functional groups was the key to success.

1753 That being said, some systematic difficulties were encountered when
1754 evaluating the efficacy of the classification models. To simplify the classification
1755 problem, we did not use all five patterns of activity identified in our grouping scheme.
1756 Instead, we used the three larger categories (Rhythmic, Weakly Rhythmic and
1757 Arrhythmic) to surmise if that was the minimum amount of labels necessary for an
1758 effective classification strategy. For both the rhythmic and arrhythmic categories, the
1759 minimum amount of information was enough to properly classify the individuals. The
1760 weakly rhythmic category, on the other hand, was consistently misclassified in every
1761 experiment. There was no perfect combination of SAX parameters and classification
1762 algorithms that would result in the consistent correct classification of the weakly
1763 rhythmic individuals. Therefore, the minimum amount of information was not
1764 sufficient to inform a proper classification model for the entirety of the *L. malachurum*
1765 dataset. Thus, in the future it would be advisable to divide the weakly rhythmic
1766 classification into smaller categories for better results, as they are probably
1767 obfuscated by the simplification of the data.

1768 The approach that we have developed for preprocessing a circadian data set
1769 before evaluation is the first of its kind to our knowledge, and as such, cannot be
1770 compared to past studies in the field. Our use of SAX for this data set is unique,
1771 since in circadian science, it is more common to use techniques like averages/rolling
1772 averages, self/auto correlations and Fourier based transformations (Refinetti et al.
1773 2007). Our results from experimenting with optimal transformation parameters using
1774 SAX were not only successful in yielding proper classifications, but also consistent
1775 with the findings in other fields, where this transformation worked best with smaller

values for its parameters (Lin and Li 2009). In the future, it would be interesting to replicate these experiments and compare the use of SAX transformations with those of wavelets, as wavelets are a commonly used tool in circadian analysis.

Furthermore, there is still the question of transferability. Other species of *Lasioglossum* should be classified using the SAX transformed data with either decision trees or KNN to confirm how generalizable is the use of our pipeline.

This work sets the basis for the use of a novel subject of study and an innovative systematic approach to preprocessing circadian data. We characterized the circadian behaviors of four never before described species of halictid bees. The results from analyzing the bees suggests that circadian behavior may have a complementary relationship with sociality. In addition, we developed and tested a unique preprocessing pipeline utilizing machine learning for the purpose of facilitating the description of organisms for whom their circadian phenotypes are unknown. Future endeavors should focus on testing the transferability of the tool. Furthermore, if we wish to strongly conclude that sociality may serve as a zeitgeber for Halictid bees, we would benefit from replications of our study with larger sample sizes and consideration of sex and seasonality, as well as a larger pool of described species.

- 1799
- 1800 **Agostinelli, F., Ceglia, N., Shahbaba, B., Sassone-Corsi, P., & Baldi, P. (2016).**
- 1801 What time is it? Deep learning approaches for circadian rhythms. *Bioinformatics*,
- 1802 32(12), i8–i17. <https://doi.org/10.1093/bioinformatics/btw243>
- 1803 **Alboukadel Kassambara and Fabian Mundt (2020).** factoextra: Extract and
- 1804 Visualize the Results of Multivariate Data Analyses. R package version 1.0.7.
- 1805 <https://CRAN.R-project.org/package=factoextra>
- 1806 **Anafi, R. C., Lee, Y., Sato, T. K., Venkataraman, A., Ramanathan, C., Kavakli, I.**
- 1807 **H., ... Hogenesch, J. B. (2014).** Machine Learning Helps Identify CHRONO as a
- 1808 Circadian Clock Component. *PLoS Biology*, 12(4), e1001840.
- 1809 <https://doi.org/10.1371/journal.pbio.1001840>
- 1810 **Beer, K., & Helfrich-Förster, C. (2020).** Post-embryonic Development of the
- 1811 Circadian Clock Seems to Correlate With Social Life Style in Bees. *Frontiers in Cell*
- 1812 *and Developmental Biology*, 8, 1325. <https://doi.org/10.3389/fcell.2020.581323>
- 1813 **Bezdek, J. C., Chuah, S. K., & Leep, D. (1986).** Generalized k-Nearest Neighbor
- 1814 rules IEEE IJCNN conference View project Big Data Cluster Analysis View project
- 1815 GENERALIZED k-NEAREST NEIGHBOR RULES*. *Fuzzy Sets and Systems*, 18,
- 1816 237–256. [https://doi.org/10.1016/0165-0114\(86\)90004-7](https://doi.org/10.1016/0165-0114(86)90004-7)
- 1817 **Bhatia, N., & Vandana. (2010).** Survey of Nearest Neighbor Techniques. *IJCSIS*)
- 1818 *International Journal of Computer Science and Information Security*, 8(2). Retrieved
- 1819 from <http://arxiv.org/abs/1007.0085>

1820 **Bloch, G., & Grozinger, C. M. (2011).** Social molecular pathways and the evolution
 1821 of bee societies. *Philosophical Transactions of the Royal Society B: Biological*
 1822 *Sciences*, Vol. 366. <https://doi.org/10.1098/rstb.2010.0346>

1823 **Bloch, G., Toma, D. P., & Robinson, G. E. (2001).** Behavioral rhythmicity, age,
 1824 division of labor and period expression in the honey bee brain. *Journal of Biological*
 1825 *Rhythms*, 16(5), 444–456. <https://doi.org/10.1177/074873001129002123>

1826 **Breiman, L. (2001)** Random Forests. *Machine Learning* 45, 5–32 .
 1827 <https://doi.org/10.1023/A:1010933404324>

1828 **Burkov, A. (2019).** The hundred-page machine learning book. Canada: Andriy
 1829 Burkov. Fumo, D. (2017, August 17).

1830 **Cordero-Martínez, C.S., et al. (2017).** The role of circadian rhythms, humidity and
 1831 temperature oscillations on the foraging patterns of *S. curvicornis* and *L.*
 1832 *malachurum*. Poster

1833 **Danforth, B. N., Conway, L., & Ji, A. S. (2003).** Phylogeny of Eusocial
 1834 *Lasioglossum* Reveals Multiple Losses of Eusociality within a Primitively Eusocial
 1835 Clade of Bees (Hymenoptera: Halictidae). *Syst. Biol*, 52(1), shrimp.
 1836 <https://doi.org/10.1080/10635150390132687>

1837 **Danforth, B. N., Eardley, C., Packer, L., Walker, K., Pauly, A., &**
 1838 **Randrianambinintsoa, F. J. (2008).** Phylogeny of Halictidae with an emphasis on
 1839 endemic African Halictinae. *Apidologie*, Vol. 39, pp. 86–101.
 1840 <https://doi.org/10.1051/apido:2008002>

1841 **Delphia, C. M., & Gibbs, J. (2020).** New Island Records for *Lasioglossum*
 1842 (Hymenoptera: Halictidae) from the Virgin Islands, West Indies. *Journal of the*
 1843 *Kansas Entomological Society*, 92(2). <https://doi.org/10.2317/0022-8567-92.2.479>
 1844 **Dubowy, C., & Sehgal, A. (2017).** Circadian rhythms and sleep in *Drosophila*
 1845 *melanogaster*. *Genetics*, 205(4), 1373–1397.
 1846 <https://doi.org/10.1534/genetics.115.185157>
 1847 **Dubowy, C., & Sehgal, A. (2017).** Circadian rhythms and sleep in *Drosophila*
 1848 *melanogaster*. *Genetics*, 205(4), 1373–1397.
 1849 <https://doi.org/10.1534/genetics.115.185157>
 1850 **Eickwort G.C. (1988).** Distribution patterns and biology of West Indian sweat bees
 1851 (Hymenoptera: Halictidae). In Liebherr J.K. (Ed.), *Zoogeography of Caribbean*
 1852 *Insects* (pp. 232–253). Ithaca: Cornell University Press.
 1853 **Eickwort G.C., Eickwort J.M, Gordon J., Eickwort M.A. (1996)** Solitary behavior in
 1854 a high-altitude population of the social sweat bee *Halictus rubicundus* (Hymenoptera:
 1855 Halictidae), *Behav. Ecol. Sociobiol.* 38, 227–233.
 1856 **Field, J. (1996)** Patterns of provisioning and iteroparity in a solitary halictine bee,
 1857 *Lasioglossum (Evylaeus) fratellum* (Perez), with notes on *L. (E.) calceatum* (Scop.)
 1858 and *L. (E.) villosulum* (K.). *Insect. Soc.* 43, 167–182
 1859 **Field, J., Paxton, R.J., Soro, A., Bridge, C. (2010).** Cryptic plasticity underlies a
 1860 major evolutionary transition. *Curr. Biol.* 20, 2028–2031.
 1861 **Geissmann Q, Garcia Rodriguez L, Beckwith EJ, Gilestro GF (2019)** Rethomics:
 1862 An R framework to analyse high-throughput behavioural data. *PLoS ONE* 14(1):
 1863 e0209331. <https://doi.org/10.1371/journal.pone.0209331>

1864 **Géron, A. (2017).** Hands-on machine learning with Scikit-Learn and TensorFlow :
 1865 concepts, tools, and techniques to build intelligent systems. In O'Reilly Media.

1866 **Giannoni-Guzmán, M. A., Aleman-Rios, J., Melendez Moreno, A. M., Diaz**
 1867 **Hernandez, G., Perez, M., Loubriel, D., ... Agosto-Rivera, J. L. (2020).** The Role
 1868 of Temperature on the Development of Circadian Rhythms in Honey Bee Workers
 1869 Authors. BioRxiv, 2020.08.17.254557. <https://doi.org/10.1101/2020.08.17.254557>

1870 **Giannoni-Guzmán, M. A., Avalos, A., Perez, J. M., Loperena, E. J. O., Kayým,**
 1871 **M., Medina, J. A., ... Agosto-Rivera, J. L. (2014).** Measuring individual locomotor
 1872 rhythms in honey bees, paper wasps and other similar-sized insects. Journal of
 1873 Experimental Biology, 217(8), 1307–1315. <https://doi.org/10.1242/jeb.096180>

1874 **Gibbs, J. (2018).** Bees of the genus *Lasioglossum* (Hymenoptera: Halictidae) from
 1875 Greater Puerto Rico, West Indies. European Journal of Taxonomy, 0(400). Retrieved
 1876 from <https://europeanjournaloftaxonomy.eu/index.php/ejt/article/view/528/1183>

1877 **Gibbs, J., Brady, S. G., Kanda, K., & Danforth, B. N. (2012).** Phylogeny of
 1878 halictine bees supports a shared origin of eusociality for *Halictus* and *Lasioglossum*
 1879 (*Apoidea*: *Anthophila*: *Halictidae*). Molecular Phylogenetics and Evolution, 65(3).
 1880 <https://doi.org/10.1016/j.ympev.2012.08.013>

1881 **Goldin, D. Q., & Kanellakis, P. C. (1995).** On similarity queries for time-series data:
 1882 Constraint specification and implementation. Lecture Notes in Computer Science
 1883 (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in
 1884 Bioinformatics), 976, 137–153. https://doi.org/10.1007/3-540-60299-2_9

1885 **Gonzalez, V. H., Pascual, C., Burrows, S., Çakmak, I., & Barthell, J. F. (2014).**
 1886 Pollen collecting behavior of *Systropha planidens* Giraud, 1861 (Hymenoptera:

1887 Halictidae) in Turkey. Pan-Pacific Entomologist, 90(4), 226–230.

1888 <https://doi.org/10.3956/2014-90.4.226>

1889 **Grozdanić, S. & Z. Vasić. (1968).** Observations on the solitary bees *Systropha*
1890 *planidens* Gir., and *Systropha curvicornis* Scop. Bulletin de Museum d’Histoire
1891 Naturelle de Belgrade 21b:133–138 [In Serbo-Croatian, with German summary].

1892 **Hans W. Borchers (2019).** pracma: Practical Numerical Math Functions. R package
1893 version 2.2.9. <https://CRAN.R-project.org/package=pracma>

1894 **Helm, B., & Visser, M. E. (2010).** Heritable circadian period length in a wild bird
1895 population. Proceedings of the Royal Society B: Biological Sciences, 277(1698).
1896 <https://doi.org/10.1098/rspb.2010.0871>

1897 **Ingram, K. K., Kutowoi, A., Wurm, Y., Shoemaker, D., Meier, R., & Bloch, G.**
1898 **(2012).** The Molecular Clockwork of the Fire Ant *Solenopsis invicta*. PLoS ONE,
1899 7(11), e45715. <https://doi.org/10.1371/journal.pone.0045715>

1900 **Jeanson, R., Clark, R., Holbrook, C., Bertram, S., Fewell, J., & Kukuk, P. (2008).**
1901 Division of labour and socially induced changes in response thresholds in
1902 associations of solitary halictine bees.
1903 <https://doi.org/10.1016/j.anbehav.2008.04.007i>

1904 **Jud, C., Schmutz, I., Hampp, G., Oster, H., & Albrecht, U. (2005).** A guideline for
1905 analyzing circadian wheel-running behavior in rodents under different lighting
1906 conditions. Biological Procedures Online, 7(1), 101–116.
1907 <https://doi.org/10.1251/bpo109>

1908 **Kriegel, H.-P., Schubert, E., & Zimek, - Arthur. (2017).** The (black) art of runtime
 1909 evaluation: Are we comparing algorithms or implementations? Knowledge and
 1910 Information Systems, 52, 341–378. <https://doi.org/10.1007/s10115-016-1004-2>
 1911 **Leonard Kaufman Peter J. Rousseeuw.(1990).** Partitioning Around Medoids
 1912 (Program PAM).Wiley Series in Probability and Statistics.
 1913 <https://doi.org/10.1002/9780470316801.ch2>
 1914 **Levine, J. D., Funes, P., Dowse, H. B., & Hall, J. C. (2002).** Signal analysis of
 1915 behavioral and molecular cycles. BMC Neuroscience, 3.
 1916 <https://doi.org/10.1186/1471-2202-3-1>
 1917 **Lin, J., & Li, Y. (2009).** Finding Structural Similarity in Time Series Data Using Bag-
 1918 of-Patterns Representation. In LNCS (Vol. 5566). Retrieved from
 1919 [https://pdfs.semanticscholar.org/5b5e/e84ac8ab484c33b3b152c9af5a0715217e53.p](https://pdfs.semanticscholar.org/5b5e/e84ac8ab484c33b3b152c9af5a0715217e53.pdf)
 1920 [df](https://pdfs.semanticscholar.org/5b5e/e84ac8ab484c33b3b152c9af5a0715217e53.pdf)
 1921 **Maechler, M., Rousseeuw, P., Struyf, A., Hubert, M., Hornik, K.(2019).** cluster:
 1922 Cluster Analysis Basics and Extensions. R package version 2.1.0.
 1923 **Michener, C. D. (1969).** Comparative Social Behavior of Bees. Annual Review of
 1924 Entomology, 14(1), 299–342. <https://doi.org/10.1146/annurev.en.14.010169.001503>
 1925 **Mistlberger, R. E., & Skene, D. J. (2004).** Social influences on mammalian
 1926 circadian rhythms: Animal and human studies. Biological Reviews of the Cambridge
 1927 Philosophical Society, Vol. 79, pp. 533–556.
 1928 <https://doi.org/10.1017/S1464793103006353>

1929 **Monti, S. (2003).** Consensus Clustering: A Resampling-Based Method for Class
 1930 Discovery and Visualization of Gene Expression Microarray Data. *Machine Learning*,
 1931 52(1/2), 91-118. doi:10.1023/a:1023949509487

1932 **Moore, D., Angel, J. E., Cheeseman, I. M., Fahrbach, S. E., & Robinson, G. E.**
 1933 **(1998).** Timekeeping in the honey bee colony: Integration of circadian rhythms and
 1934 division of labor. *Behavioral Ecology and Sociobiology*, 43(3), 147–160.
 1935 <https://doi.org/10.1007/s002650050476>

1936 **Ortiz-Alvarado, Carlos A., Ambrose, Alexandria F., Cordero-Martinez, Claudia**
 1937 **S., Silva-Echeandia, Sebastian A., Petanidou, Theodora, Tscheulin Thomas,**
 1938 **Gonzalez, Victor H., Giray, Tugrul, Agosto-Rivera, Jose A., Barthell, John F. (in**
 1939 **rev)** Daily foraging activity and circadian rhythms of a guild of three large carpenter
 1940 bee species in Chasteberry (*Vitex agnus-castus* L.) on a Greek island (Lesvos).

1941 **Patiny, S., & Michez, D. (2007).** New insights on the distribution and floral choices
 1942 of *Systropha* Illiger, 1806 in Africa (Hymenoptera, Apoidea), with description of a
 1943 new species from Sudan. *Zootaxa*, (1461), 59–68.
 1944 <https://doi.org/10.11646/zootaxa.1461.1.6>

1945 **Patiny, S., Michez, D., & Danforth, B. N. (2008).** Phylogenetic relationships and
 1946 host-plant evolution within the basal clade of Halictidae (Hymenoptera, Apoidea).
 1947 *Cladistics*, 24(3), 255–269. <https://doi.org/10.1111/j.1096-0031.2007.00182.x>

1948 **Peters, R. S., Krogmann, L., Mayer, C., Donath, A., Gunkel, S., Meusemann, K.,**
 1949 **... Niehuis, O. (2017).** Evolutionary History of the Hymenoptera. *Current Biology*,
 1950 27(7). <https://doi.org/10.1016/j.cub.2017.01.027>

1951 **Radmacher, S., & Strohm, E. (2010).** Factors affecting offspring body size in the
 1952 solitary bee *Osmia bicornis* (Hymenoptera, Megachilidae) Factors affecting offspring
 1953 body size in the solitary bee *Osmia bicornis* (Hymenoptera, Megachilidae)*.
 1954 *Apidologie*, 41, 169–177. <https://doi.org/10.1051/apido/2009064>

1955 **R Core Team (2020).** R: A language and environment for statistical computing. R
 1956 Foundation for Statistical Computing, Vienna, Austria. URL [https://www.R-](https://www.R-project.org/)
 1957 [project.org/](https://www.R-project.org/).

1958 **Refinetti, Roberto & Lissen, Germaine & Halberg, Franz. (2007).** Procedures for
 1959 numerical analysis of circadian rhythms. *Biological rhythm research*. 38. 275-325.
 1960 10.1080/09291010600903692.

1961 **Richards, M. H. (2000).** Evidence for geographic variation in colony social
 1962 organization in an obligately social sweat bee, *Lasioglossum malachurum* Kirby
 1963 (Hymenoptera; Halictidae). *Canadian Journal of Zoology*, 78(7).
 1964 <https://doi.org/10.1139/z00-064>

1965 **Richards, M. H., von Wettberg, E. J. & Rutgers, A. C. (2003).** A novel social
 1966 polymorphism in a primitively eusocial bee. *Proc. Natl Acad. Sci. USA* 100, 7175–
 1967 7180. doi:10.1073/pnas.1030738100

1968 **Roenneberg, T., Daan, S., and Mellow, M. (2003).** The art of entrainment. *J. Biol.*
 1969 *Rhythms* 18, 183–194.

1970 **Roenneberg, T., Wirz-Justice, A., & Mellow, M. (2003).** Life between Clocks: Daily
 1971 Temporal Patterns of Human Chronotypes. *Journal of Biological Rhythms*, 18(1).
 1972 <https://doi.org/10.1177/0748730402239679>

1973 **Rubin, E. B., Shemesh, Y., Cohen, M., Elgavish, S., Robertson, H. M., & Bloch,**
 1974 **G. (2006).** Molecular and phylogenetic analyses reveal mammalian-like clockwork in
 1975 the honey bee (*Apis mellifera*) and shed new light on the molecular evolution of the
 1976 circadian clock. *Genome Research*, 16(11), 1352–1365.
 1977 <https://doi.org/10.1101/gr.5094806>

1978 **Schwarz, M. P., Richards, M. H., & Danforth, B. N. (2006).** Changing Paradigms in
 1979 Insect Social Evolution: Insights from Halictine and Allodapine Bees.
 1980 <https://doi.org/10.1146/annurev.ento.51.110104.150950>

1981 **Schwarz, M. P., Richards, M. H., & Danforth, B. N. (2006).** Changing Paradigms in
 1982 Insect Social Evolution: Insights from Halictine and Allodapine Bees.
 1983 <https://doi.org/10.1146/annurev.ento.51.110104.150950>

1984 **Steitz, I., Kingwell, C., Paxton, R. J., & Ayasse, M. (2018).** Evolution of Caste-
 1985 Specific Chemical Profiles in Halictid Bees. *Journal of Chemical Ecology*, 44(9),
 1986 827–837. <https://doi.org/10.1007/s10886-018-0991-8>

1987 **Toth, A. L., & Rehan, S. M. (2017).** Molecular Evolution of Insect Sociality: An Eco-
 1988 Evo-Devo Perspective. *Annual Review of Entomology*, 62(1).
 1989 <https://doi.org/10.1146/annurev-ento-031616-035601>

1990 **VanderPlas, J. T. (2017).** Understanding the lomb-scargle periodogram. *ArXiv*, Vol.
 1991 236, p. 16. <https://doi.org/10.3847/1538-4365/aab766>

1992 **West-Eberhard MJ. (2003).** Developmental Plasticity and Evolution. Oxford, UK:
 1993 Oxford Univ. Press

1994 **Westrich, P. (1989).** Die Wildbienen Baden-Wu"rttembergs. Ulmer Verlag, Stuttgart,
 1995 Germany. 972 pp

1996 **Wickham, H. (2009).** ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag
 1997 New York. In Media (Vol. 35).

1998 **Wilkerson, D. M, Hayes, N. D. (2010).** "ConsensusClusterPlus: a class discovery
 1999 tool with confidence assessments and item tracking." Bioinformatics, 26(12), 1572-
 2000 1573. <http://bioinformatics.oxfordjournals.org/content/26/12/1572.abstract>.

2001 **Wyman, L. M., & Richards, M. H. (2003).** Colony social organization of
 2002 *Lasioglossum malachurum* Kirby (Hymenoptera, Halictidae) in southern Greece.
 2003 *Insectes Sociaux*, 50(3). <https://doi.org/10.1007/s00040-003-0647-7>

2004 **Yuan Tang, Masaaki Horikoshi, and Wenxuan Li.(2016).** "ggfortify: Unified
 2005 Interface to Visualize Statistical Result of Popular R Packages." *The R Journal* 8.2 :
 2006 478-489.

2007

2008

2009

2010

2011

2012

2013

2014

2015

2016

2017

2018

2019

2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037

Appendix

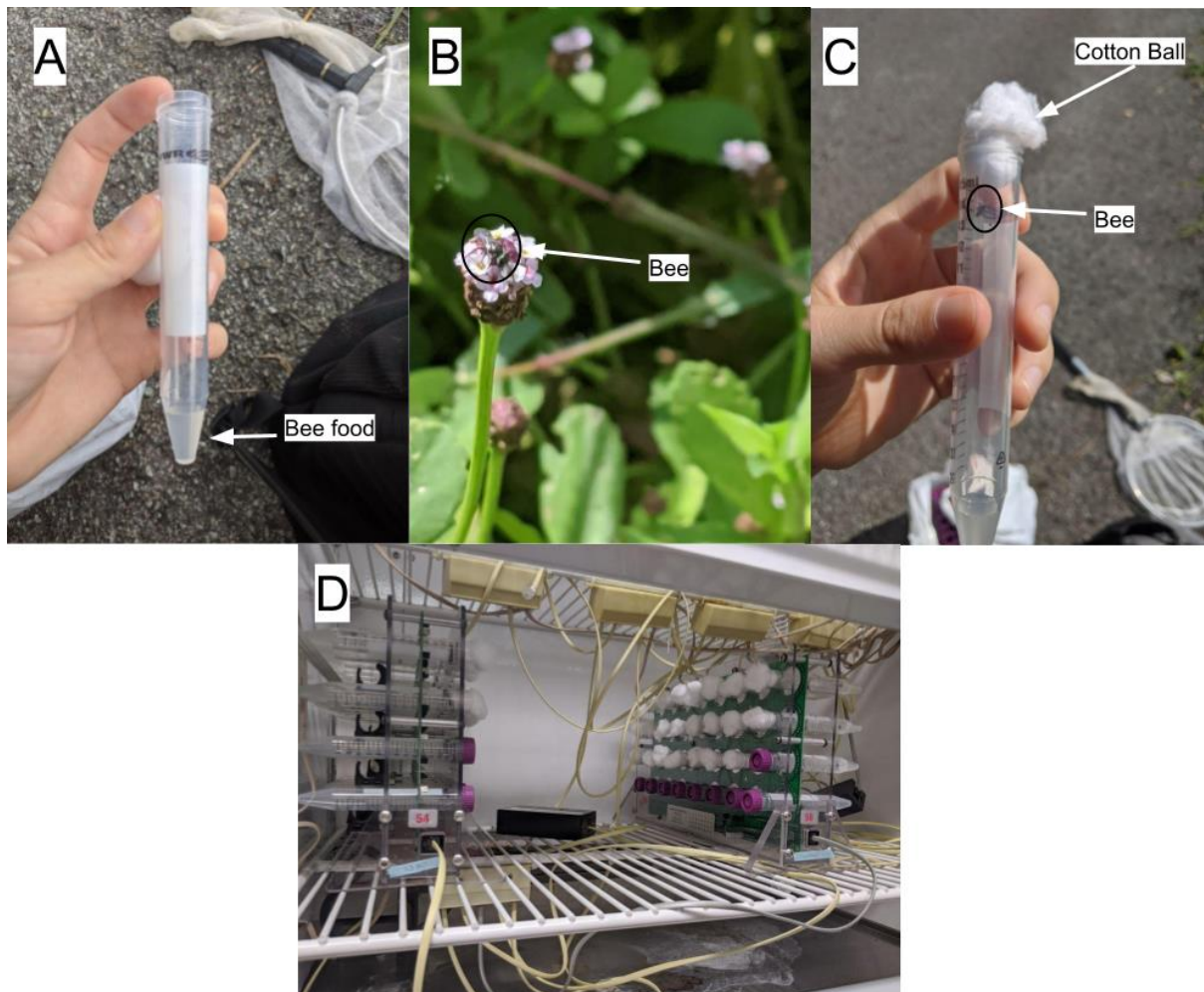


Figure A.1: Methods for capturing and husbandry of the bees. A.1.A) Shows a tube that will house an individual bee, with the agar/sucrose food gel. A.1.B) A bee specimen visiting *Phyla nodiflora*, where they will be captured as illustrated in A.1.C) with cotton instead of a cap. A.1.D) The bees in their final destination in the incubator placed in their monitors.