

BioLockJ User Manual



College of Computing and Informatics
Department of Bioinformatics
Fodor Lab

Change Log

#	Date	Description
1.0	07/27/2017	Initial document BioLockJ v1.0

File Name

BioLockJ_User_Manual.pdf

Table of Contents

1. Pipeline Overview	3
1.1 System Diagram.....	4
1.2 Installation.....	4
1.3 Software Dependencies	4
1.4 Launching BioLockJ.....	6
1.5 Failure Recovery	6
1.6 Email Notification	7
2. Input File Specification	7
2.1 Sequence Files	7
2.2 Metadata	8
2.3 Descriptor	8
2.4 BioLockJ Configuration	8
2.4.1 <i>Example BioLockJ Configuration</i>	13
3. System Architecture	16
3.1 Java Project.....	16
3.1.1 <i>Bundled JAR Files</i>	17
3.1.2 <i>Logging Framework</i>	18
3.1.4 <i>Log Levels</i>	18
3.2 BioLockJ Module Overview	19
3.2.1 <i>Module Directory Structure</i>	20
3.2.2 <i>Module Output</i>	21
4. QIIME Modules	22
4.1 Qiime Preprocessor	22
4.2 Closed Reference Classifier	23
4.3 Merge OTU Tables	24
4.4 Open Reference Classifier	25
4.5 De Novo Classifier.....	27
4.6 QIIME Classifier	29
5. Extending the Pipeline	30
5.1 Enhancing the R Script.....	30
5.2 Comparing Classifiers	30
5.3 Adding New Classifiers	30
Appendix A: Contact Information.....	31

1. Pipeline Overview

BioLockJ is a light-weight, extensible, metagenomics pipeline designed to improve the speed, accuracy, and reproducibility of 16s amplicon and whole genome sequencing (WGS) data analysis. BioLockJ runs on any Linux system (and by extension OSX) but is most powerful in a high performance computing environment by utilizing its parallel processing capabilities.

Pipeline execution is guided by a single BioLockJ configuration file which can be used to reproduce your analysis and serves to document all runtime parameters.

Note: BioLockJ properties listed in [Section 2.4](#) will appear in italics throughout this text.

Primary Inputs

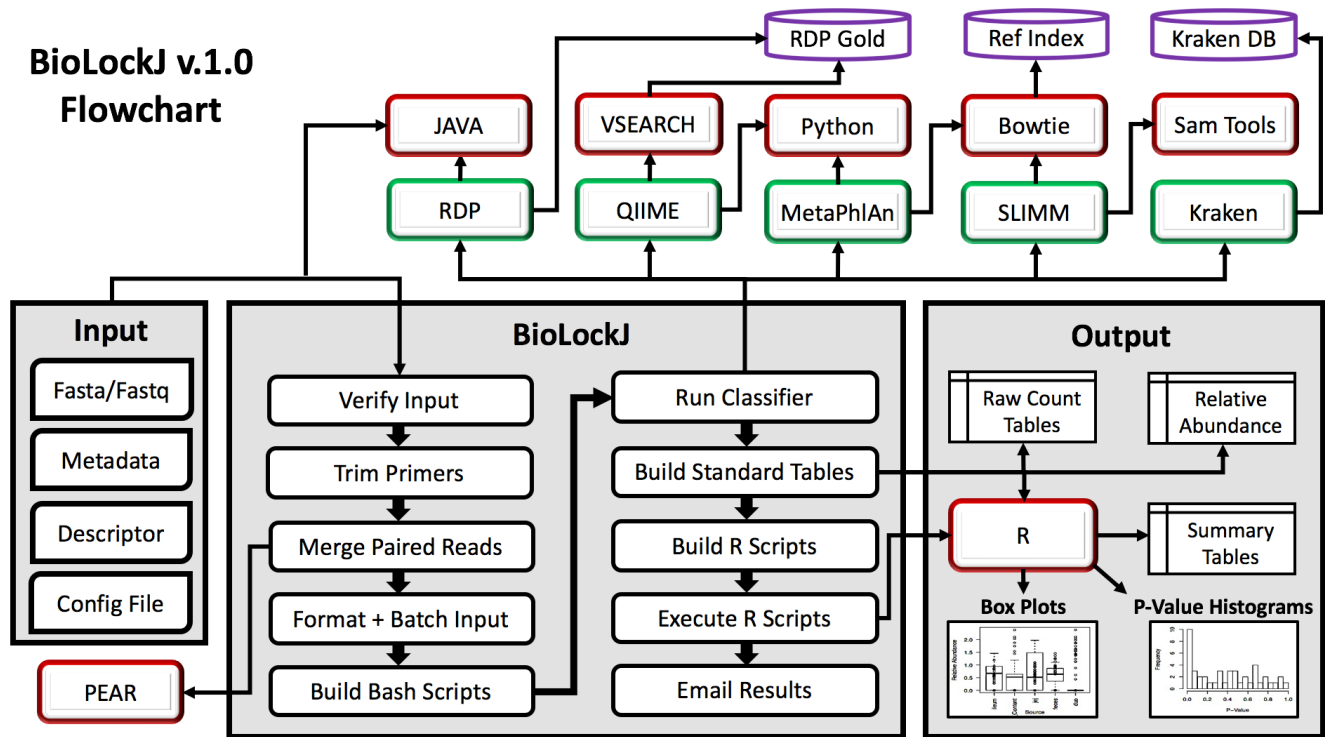
- | | |
|-------------------------------|------------------------------|
| 1. <i>input.dirs</i> | Fast-A/Fast-Q sequence files |
| 2. <i>metadata.file</i> | Path to metadata file |
| 3. <i>metadata.descriptor</i> | Path to descriptor file |

BioLockJ writes and executes bash shell scripts that call command line bioinformatics tools based on user input, provided via the BioLockJ configuration file. Generated scripts are reliable, organized, efficient, and reproducible.

Execution Summary

1. Formats sequences to meet classifier specifications
2. Trim primers, merge paired reads, and rarefy as needed
3. Classifies sequences using RDP, QIIME, Kraken, MetaPhlAn, or SLIMM
4. Generates raw count and relative abundance tables
5. Builds statistical models to find significant OTUs correlated with metadata
6. Generates summary tables and PDF reports with histograms and bar charts
7. Emails a summary report upon completion

1.1 System Diagram



1.2 Installation

- BioLockJ is deployed as a single JAR file, which can be executed by JAVA
 1. Download GitHub BioLockJ project from github.com/mikesioda/BioLockJ
 2. Install software dependencies

1.3 Software Dependencies

- Required software: Java, R, R packages (Kendall & Coin)
- BioLockJ can be deployed with as few as one classifier, the rest are optional
- Unused classifiers and their dependencies do not need to be installed

BioLockJ Software Dependencies

#	Program	Version	Description
1	Java	1.8	Required - java.com
2	Python	2.7.12	Required by QIIME and MetaPhlAn - python.org
3	Bowtie2	2.3.2	Required by MetaPhlAn and SLIMM bowtie-bio.sourceforge.net/bowtie2/index.shtml
4	SAMtools	1.4	Required by SLIMM - htslib.org
5	Vsearch	2.4.3	Required for chimera removal in QIIME. github.com/torognes/vsearch
6	PEAR	0.9.8	Paired-End reAd merger, used by RDP & QIIME sco.h-its.org/exelixis/web/software/pear
7	RDP	2.11	16S Classifier: Ribosomal Database Project github.com/rdpstaff/classifier
8	QIIME	1.9.1	16S Classifier: Quantitative Insights Into Microbial Ecology qiime.org
9	Open MPI	1.10	Open Message Passing Interface, required by QIIME open-mpi.org
10	Kraken	0.10.5-beta	WGS Classifier - ccb.jhu.edu/software/kraken
11	MetaPhlAn2	2.0	WGS Classifier: Metagenomic Phylogenetic Analysis huttenhower.sph.harvard.edu/metaphlan2
12	SLIMM	0.2.2	WGS Classifier: Species Level Identification of Microbes from Metagenomes - github.com/seqan/slimm
13	GNU Awk	4.0.2	Convert Fastq to Fasta for QIIME - gnu.org/software/gawk
14	GNU Gzip	1.5	Decompress gzipped files for QIIME - gnu.org/software/gzip
15	R	3.2.3	Statistical modeling package - cran.r-project.org
16	Kendall	2.2	Kendall rank correlation p-values for continuous data types cran.r-project.org/web/packages/Kendall/index.html
17	Coin	1.2	Conditional Inference Procedures in a Permutation Test Framework: Computes exact p-value for Wilcox_test cran.r-project.org/web/packages/coin/index.html

1.4 Launching BioLockJ

- Runtime parameters are located in the BioLockJ configuration file
- To run BioLockJ, call java on BioLockJ.jar and pass the configuration file path

Java Command

```
nohup java -jar ${jar_path}/BioLockJ.jar ${config_file_path} > /dev/null 2>&1 &
```

Command Part	Description
nohup	Continue execution after the terminal connection is closed
java	Java executable command
-jar	Required parameter to run a JAR
<code>\${jar_path}/BioLockJ.jar</code>	Path to BioLockJ.jar
<code>\${config_file_path}</code>	Path to configuration file
<code>> /dev/null 2>&1</code>	Discard terminal output – it is already output to \${LOG_FILE}
<code>&</code>	Run process in background to free terminal for other work

1.5 Failure Recovery

- Save time by recovering from failures rather than restarting the entire pipeline
- To determine root cause of failure:
 - Check log file for Java stack trace pointing to line number in Java class
 - Check failed module's "failures" directory for failure indicator files
- Fix the problem and resume pipeline by updating BioLockJ configuration file:
 - Set successfully completed module control property values = N
 - Set input.dirs = output directory of last successful module
 - Set metadata.file = path of current version if updated by pipeline
 - Set metadata.descriptor = path of current version if updated by pipeline
- Launch BioLockJ with the updated configuration

1.6 Email Notification

- BioLockJ sends *email.to* recipients a summary report & log file after pipeline execution
- Email body reports runtime for each module
- Email body reports script failures, if any
- Log4J attachment contains execution details based on log level (see [Section 3.3.2](#))
- Option to attach Qsub output and error files if run in clustered environment



2. Input File Specification

2.1 Sequence Files

- Supported classifiers will accept 2 sequence file formats:
 1. Fast-A
 2. Fast-Q
- Files may be gzipped (with extension .gz)
- Files names must be unique – important to consider with multiple *input.dirs*
- Kraken & RDP sequence files may be multiplexed

2.2 Metadata

- Columns must be tab-delimited
- First column must hold the Sample ID (regardless of column header name)
- Blank cells are considered empty Strings
- Null values must be indicated using *metadata.nullValue*

2.3 Descriptor

- Columns must be tab-delimited
 - Attribute (Column 1): required field, contains name of metadata attribute
 - Type (Column 2): required field, options: Binary, Continuous, Categorical
 - Comments (Column 3): optional field, for user comments
- Metadata first column contains Sample ID, not to be included in descriptor file
- All metadata columns, other than the 1st, must be assigned an attribute type
- Type determines the statistical model implemented in R for each *report.attributes*:
 - Binary: Exact Wilcoxon Signed Rank Sum (Coin package)
 - Continuous: Kendall Tau Rank Correlation (Kendall package)
 - Categorical: One-Way ANOVA

Attribute	Type	Comments
attribute1	Binary	Metadata column must contain 1 – 2 unique values
attribute2	Continuous	Metadata column must contain numeric values
attribute3	Categorical	Metadata column must contain 2 or more unique values

2.4 BioLockJ Configuration

- The BioLockJ configuration file contains runtime parameters as name-value pairs:
- List properties are comma separated

```
property1.name=property.value1
property2.name=property.value2a,property.value2b,property.value2c
...
```


Property	Value
<i>project.name</i>	Each pipeline execution generates a timestamped directory: project.name_yyyyMMdd_kkmmss
<i>project.rootDir</i>	Parent directory for BioLockJ project run directories
<i>project.copyInputFiles</i>	Options: Y/N. If Y, copy <i>input.dirs</i> into the project directory
<i>project.deleteTempFiles</i>	Options: Y/N. If Y, delete module temp dirs after execution
<i>project.classifierType</i>	Options: rdp, qiime, kraken, metaphlan, slimm
<i>control.runOnCluster</i>	Options: Y/N. If Y, cluster properties is required
<i>control.trimSeqs</i>	Options: Y/N. If Y, the SeqTrimmer module will execute
<i>control.mergePairs</i>	Options: Y/N. If Y, the PairedSeqMerger module will execute
<i>control.rarefySeqs</i>	Options: Y/N. If Y, the Rarefier module will execute
<i>control.runClassifier</i>	Options: Y/N. If Y, the ClassifierModule will execute
<i>control.runParser</i>	Options: Y/N. If Y, the ParserModule will execute
<i>control.run_rScript</i>	Options: Y/N. If Y, the RScriptBuilder will execute
<i>input.dirs</i>	Must contain files expected by the first module executed. Multiple directories must be comma-separated.
<i>input.ignoreFiles</i>	Files listed here will be ignored if found in <i>input.dirs</i> . Multiple files must be comma-separated.
<i>input.demultiplex</i>	Options: Y/N. If Y, sequence files include reads from multiple samples & sample IDs must be extracted from the sequence headers. RDP & Kraken classifiers only. If N, there is one sample per file and the file name must contain the sample ID.
<i>input.pairedReads</i>	Options: Y/N. If Y, file names must include <i>input.forwardFileSuffix</i> or <i>input.reverseFileSuffix</i>
<i>input.forwardFileSuffix</i>	File name suffix to indicate a forward read
<i>input.reverseFileSuffix</i>	File name suffix to indicate a reverse read
<i>input.trimPrefix</i>	For files named by Sample ID, provide the prefix preceding the ID to trim when extracting Sample ID. If <i>input.demultiplex</i> =Y, provide any characters in the sequence header preceding the ID. For fastq, typically "@".

Property	Value
<i>input.trimSuffix</i>	For files named by Sample ID, provide the suffix after the ID, often this is just the file extension. Do not include read direction indicators listed in <i>input.forwardFileSuffix/reverseFileSuffix</i> . If <i>input.demultiplex=Y</i> , provide 1st character in the sequence header after ID; for fastq, typically ":" char
<i>input.rarefyMinNumSeqs</i>	Discard samples without min # of seqs
<i>input.rarefyMaxNumSeqs</i>	Randomly select max # of seqs for each sample
<i>input.trimSeqPath</i>	Path to file containing primers to trim File must contain only one sequence per line
<i>cluster.batchCommand</i>	The command to submit jobs on the cluster
<i>cluster.params</i>	Include in header of scripts submitted on cluster
<i>cluster.validateParams</i>	Options: Y/N. If Y, validate procs= <i>script.numThreads</i>
<i>cluster.modules</i>	List of modules to load before execution Adds "module load" command to bash scripts
<i>script.exitOnError</i>	Options: Y/N. If Y, program exits if any script failures occur, otherwise failures logged to failure directory
<i>script.batchSize</i>	Number of sequence files to process per script
<i>script.chmodCommand</i>	Command to grant script execute permissions
<i>script.numThreads</i>	Passed to number of threads parameter in classifier
<i>metadata.file</i>	Metadata file path, attributes referenced in R properties validated based on descriptor value
<i>metadata.descriptor</i>	Descriptor file path, defines all metadata columns
<i>metadata.nullValue</i>	Define how null values are represented in metadata
<i>metadata.commentChar</i>	Define how comments are indicated in metadata
<i>report.numHits</i>	Options: Y/N. If Y, and add #Hits to metadata
<i>report.numReads</i>	Options: Y/N. If Y, and add #Reads to metadata
<i>report.fullTaxonomyNames</i>	Options: Y/N. If Y, ParserModule will use full taxonomy names in output tables
<i>report.addGenusToSpeciesName</i>	Options: Y/N. If Y, ParserModule adds genus prefix to species name in output tables

Property	Value
<i>report.useGenusFirstInitial</i>	Options: Y/N. If Y, ParserModule adds genus 1 st initial prefix to species name in output tables
<i>report.attributes</i>	R script adds statistical models for metadata attributes
<i>report.minOtuCount</i>	ParserModule ignores OTU counts below min count
<i>report.emptySpaceDelim</i>	Reports separate genus and species using this value
<i>report.taxonomyLevels</i>	Options: domain, phylum, class, order, family, genus, species. Generate reports for listed taxonomy levels
<i>email.sendNotification</i>	Options: Y/N. If Y, sent notification email
<i>email.sendQsub</i>	Options: Y/N. If Y, attach qsub output and error files
<i>email.maxAttachmentSizeMB</i>	Max size (in MB) for log file attachment
<i>email.encryptedPassword</i>	<p>Encrypted password from <i>email.from</i> account. If BioLockJ is passed a 2nd parameter (in addition to the config file), the 2nd parameter should be the clear-text password. The password will be encrypted and stored in the prop file for future use.</p> <p>WARNING: Base64 encryption is only a trivial roadblock for malicious users. This functionality is intended merely to keep clear-text passwords out of the configuration files and should only be used with a disposable <i>email.from</i> account.</p>
<i>email.from</i>	Notification emails sent from this account, provided <i>email.encryptedPassword</i> is valid
<i>email.to</i>	Comma-separated email recipients list
<i>r.logNormal</i>	Options: Y/N. If Y, use relative abundance table in R
<i>r.logBase</i>	Options: 10/e. If e, use natural log (base e), otherwise use log base 10
<i>r.maxTitleSize=25</i>	Report OTU names trimmed after max # characters
<i>r.rareOtuThreshold</i>	If >1, R will filter OTUs below # provided. If <1, R will treat # as percentage and ignore OTUs not found in that percentage of table rows in each taxa level
<i>r.filterAttributes</i>	Ordered list of attributes to report on in R script

Property	Value
<i>r.filterOperators</i>	Ordered list of logical operators to apply to <i>r.filterAttributes</i>
<i>r.filterValues</i>	Ordered list of values to compare with <i>r.filterAttributes</i>
<i>r.filterNaAttributes</i>	Rows with NA values for attributes ignored in R script
<i>r.numHistogramBreaks</i>	Number of breaks in P-value histograms output by R script
<i>exe.classifier</i>	Classifier executable command
<i>exe.classifierParams</i>	Optional classifier parameters, excluding parameters generated by BioLockJ (input files, output files, #threads)
<i>exe.rScript</i>	Executable RScript command
<i>exe.java</i>	Executable Java command
<i>exe.python</i>	Executable python command
<i>exe.gzip</i>	Executable gzip command
<i>exe.awk</i>	Executable awk command
<i>exe.bowtie</i>	Executable bowtie2 command
<i>exe.bowtie_params</i>	Optional bowtie2 parameters
<i>exe.pear</i>	Executable PEAR command
<i>exe.pear_params</i>	Optional PEAR parameters
<i>exe.samtools</i>	Executable samtools command
<i>exe.vsearch</i>	Executable vsearch command
<i>exe.vsearchParams</i>	Optional vsearch parameters
<i>rdp.minThresholdScore</i>	Required RDP minimum threshold score for valid OTUs
<i>qiime.pickOtuScript</i>	Options: pick_closed_reference_otus.py, pick_de_novo_otus.py, pick_open_reference_otus.py
<i>qiime.alphaDiversityMetrics</i>	Options listed online: scikit-bio.org/docs/latest/generated/skbio.diversity.alpha.html
<i>qiime.formatMetadata</i>	Options: Y/N. If Y, format metadata to build QIIME mapping
<i>qiime.preprocessInput</i>	Options: Y/N. If Y, decompress gzipped files and/or convert fastq files into fasta format, as required by QIIME
<i>qiime.pickOtu</i>	Options: Y/N. If Y, execute <i>qiime.pickOtuScript</i>

Property	Value
<i>qiime.mergeOtuTables</i>	Options: Y/N. If Y, merge OTU tables generated by batched scripts calling <code>pick_closed_reference_otus.py</code>
<i>qiime.removeChimeras</i>	Options: Y/N. If Y, remove chimeras after open or de novo OTU picking using <code>exe.vsearch</code>
<i>kraken.db</i>	Path to Kraken database
<i>slimm.db</i>	Path to SLIMM database
<i>slimm.refGenomeIndex</i>	Path the bowtie2 reference genome index

2.4.1 Example BioLockJ Configuration

```

project.name=twinStudy_tp2
project.rootDir=/research/microbiome/biolockj
project.copyInputFiles=N
project.deleteTempFiles=N
project.classifierType=RDP

control.runOnCluster=Y
control.trimSeqs=Y
control.mergePairs=Y
control.rarefySeqs=Y
control.runClassifier=Y
control.runParser=Y
control.run_rScript=Y

input.dirs=/datasets/16s/twinStudy/fw,/datasets/16s/twinStudy/rv
input.ignoreFiles=Cleandata.stat
input.demultiplex=N
input.pairedReads=Y
input.forwardFileSuffix=_R1
input.reverseFileSuffix=_R2
input.trimPrefix=timepoint2
input.trimSuffix=.fq
input.rarefyMinNumSeqs=10000
input.rarefyMaxNumSeqs=100000
input.trimSeqPath=/research/microbiome/primers/twinStudyPrimers.txt

```

```
cluster.batchCommand=qsub -q copperhead
cluster.params=#PBS -l procs=8,mem=32GB,walltime=12:00:00
cluster.validateParams=Y
cluster.modules=rdp/2.12
script.exitOnError=Y
script.batchSize=8
script.chmodCommand=chmod 774
script.numThreads=8

metadata.file=/research/microbiome/metadata/twinStudyMetadata.txt
metadata.descriptor=/research/microbiome/twinStudyDescriptor.txt
metadata.nullValue=NA
metadata.commentChar=##

report.numHits=Y
report.numReads=Y
report.fullTaxonomyNames=N
report.addGenusToSpeciesName=N
report.useGenusFirstInitial=Y
report.attributes=maritalStatus, sex, bmi, age
report.minOtuCount=2
report.emptySpaceDelim=.
report.taxonomyLevels=phylum,class,order,family,genus

email.sendNotification=Y
email.sendQsub=N
email.maxAttachmentSizeMB=5
email.encryptedPassword=SlrotqvCPGsFhWkKxtpwkQ==
email.from=biolockj@gmail.com
email.to=msioda@uncc.edu

r.logNormal=Y
r.logBase=e
r.maxTitleSize=25
r.rareOtuThreshold=0.25
r.filterAttributes=age
r.filterOperators=<
r.filterValues=17
r.filterNaAttributes=maritalStatus
r.numHistogramBreaks=20

exe.classifier=/apps/rdp_2.12/dist/classifier.jar
```

```
exe.classifierParams=/databases/silva128/rRNAClassifier.properties
exe.rScript=/apps/pkg/R-3.2.3/rhel7_u2-x86_64/gnu/bin/Rscript
exe.java=java
exe.python=python
exe.gzip=gzip
exe.awk=awk
exe.bowtie=bowtie2
exe.bowtie_params=no-unal, k 60
exe.pear=/apps/pear/pear-0.9.10-bin-64
exe.pear_params=t 150
exe.samtools=samtools
exe.vsearch=/apps/vsearch-2.4.3-linux-x86_64/bin/vsearch
exe.vsearchParams=db /databases/rdp_gold.fa

rdp.minThresholdScore=50

qiime.pickOtuScript=pick_closed_reference_otus.py
qiime.alphaDiversityMetrics=shannon,chao1,observed_species
qiime.formatMetadata=Y
qiime.preprocessInput=Y
qiime.pickOtus=Y
qiime.mergeOtuTables=Y
qiime.removeChimeras=N

kraken.db=/databases/kraken/all_bacteria_archaea_20170502

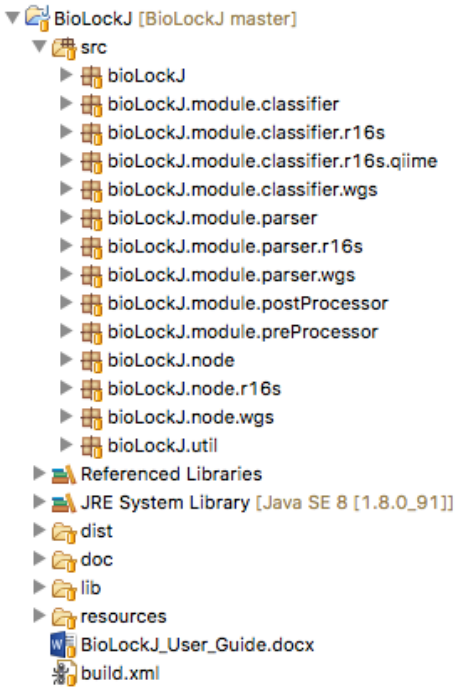
slimm.db=/apps/slimm/slimmDB_13K
slimm.refGenomeIndex=/databases/slimm/AB_13K_ref_genomes_bowtie2/AB_13K
```

3. System Architecture

- BioLockJ's modular design serves multiple purposes:
 1. Organize analysis by creating a logical separation of tasks and output file
 2. Facilitate failure recovery, easily identify where to restart pipeline after failure
 3. Promote pipeline extension, add modules to support new classifiers,

3.1 Java Project

- Screenshot from Eclipse Package Explorer (Java IDE)

	BioLockJ	BioLockJ root directory
	src	Java source code
	dist	BioLockJ.jar ANT build script target
	doc	Javadoc directory
	lib	Required Java libraries (JAR files)
	resources	<ul style="list-style-type: none"> ➤ Config file templates ➤ Metadata files ➤ Descriptor files ➤ Primer files ➤ log4j.properties
	BioLockJ User Manual	
	build.xml	Apache ANT build script

3.1.1 Bundled JAR Files

- JAR files are located in the lib directory

JAR File	Java Classes	BioLockJ Invoking Method
commons-configuration-1.10	PropertiesConfiguration PropertiesConfigurationLayout	MailUtil.encryptAndStoreEmailPassword
commons-csv-1.4	CSVFormat CSVParser CSVRecord	MetadataUtil.processFile
commons-io-2.5	FileUtils IOFileFilter NameFileFilter TrueFileFilter WildcardFileFilter	Module.initInputFiles Module.setModuleInput ApplicationManager.copyInputDirs MailUtil.getAttachments
commons-lang-2.6	NestableException	Module.initInputFiles Module.setModuleInput ApplicationManager.copyInputDirs MailUtil.getAttachments
javax.mail	BodyPart InternetAddress Message MimeBodyPart MimeMessage MimeMultipart Multipart PasswordAuthentication Session Transport	MailUtil.sendEmailNotification
log4j-1.2.17 slf4j-api-1.7.22 slf4j-log4j12-1.7.22	Logger LoggerFactory	BioLockJ.initializeGlobalProps

3.1.2 Logging Framework

- Simple Logging Facade for Java (SLF4J) wraps the chosen logging framework
- The logging framework implemented for BioLockJ is Log4J
- Adjust log specificity by setting the log level in log4j.properties (default = INFO)
- System property `{LOG_FILE}` is set based on configuration property: *project.name*
- Logger initialized in `bioLockJ.ApplicationManager.buildNewProject` method

3.1.3 Log4J Configuration

```
# Set log level and targets
log4j.rootLogger=INFO, file, stdout

# Configure command line output
log4j.appender.stdout=org.apache.log4j.ConsoleAppender
log4j.appender.stdout.Target=System.out
log4j.appender.stdout.layout=org.apache.log4j.PatternLayout
log4j.appender.stdout.layout.ConversionPattern=%d %-5p - %m%n

# Configure log file output
log4j.appender.file=org.apache.log4j.FileAppender
log4j.appender.file.File=${LOG_FILE}
log4j.appender.file.Append=false
log4j.appender.file.layout=org.apache.log4j.PatternLayout
log4j.appender.file.layout.ConversionPattern=%d %-5p - %m%n
```

3.1.4 Log Levels

Level	Log Level Description
DEBUG	All messages are logged, may result in large log files & impact performance
INFO	Informational, warning, and error messages are logged
WARN	Warning and error messages are logged
ERROR	Only error messages are logged

3.2 BioLockJ Module Overview

- Module execution occurs one at a time and in the order specified in the table below
- A module will run if it's control property is set to Y
- The first module configured to run uses *input.dirs* as input
- Additional modules use output from the previous module as input
- BioLockJ modules extend *bioLockJ.Module.java*

#	Module	Description	Control Property
1	SeqTrimmer	Remove primers from sequence files	<i>control.trimSeqs</i>
2	PairedSeqMerger	Merge paired reads	<i>control.mergePairs</i>
3	Rarefier	Discard samples with < min # seqs Randomly select max # seqs/sample	<i>control.rarefySeqs</i>
4	QiimePreProcessor	Decompress gzipped files Convert fastq to fasta Generate QIIME mapping file	<i>qiime.preprocessInput</i>
5a.1	ClosedRefClassifier	Batch reads for parallel processing Pick closed ref OTUs for QIIME	<i>qiime.pickOtu</i>
5a.2	MergeOtuTables	Merge OTU tables from 5a.1	<i>qiime.mergeOtuTables</i>
5b	OpenRefClassifier	Pick open ref OTUs for QIIME	<i>qiime.pickOtu</i>
5c	DeNovoClassifier	Pick de novo OTUs for QIIME	<i>qiime.pickOtu</i>
6	ClassifierModule	Build raw count and relative abundance tables	<i>control.runClassifier</i>
7	ParserModule	Merge metadata with classifier output to generate raw count and relative abundance tables for each level in: <i>report.taxonomyLevels</i>	<i>control.runParser</i>
8	RScriptBuilder	Build summary tables, p-value histograms, box-plots in R	<i>control.run_rScript</i>

**** Modules 4 – 5 are QIIME specific ****

3.2.1 Module Directory Structure

--- <i>project.rootDir</i>	Defined in configuration file
--- <i>project.name_<ts></i>	Defined in configuration file & <yyyyMMdd_kkmmss>
--- #_module.name	# module (in the run order) & module name
--- failures	Contains empty files indicating cause of script failures
--- output	Module output (also input for next module)
--- qsub	Contains output/error logs for each script run on cluster
--- scripts	Contains module scripts & empty status indicator files
--- temp	Contains intermediate output not passed next module

Directory	Usage
failures	Exists if module runs bash scripts
output	Exists for all modules
qsub	Exists if module bash scripts run on cluster
scripts	Exists if module runs bash or R scripts
temp	Exists if module stores intermediate output not passed to next module

Example Module Layout	Control Property
--- projects	<i>project.rootDir=/projects</i>
--- twinStudy_20170629_104521	<i>project.name=twinStudy</i>
--- 0_SeqTrimmer	<i>control.trimSeqs=Y</i>
--- output	
--- 1_PairedSeqMerger	<i>control.mergePairs=Y</i>
--- failures	
--- output	
--- temp	
--- qsub	<i>control.runOnCluster=Y</i>
--- scripts	
--- 2_Rarefier	<i>control.rarefySeqs=Y</i>
--- output	
--- temp	
--- 3_RdpClassifier	<i>control.runClassifier=Y</i>
--- failures	
--- output	
--- qsub	<i>control.runOnCluster=Y</i>
--- scripts	
--- 4_RdpParser	<i>control.runParser=Y</i>
--- output	
--- temp	
--- 5_RScriptBuilder	<i>control.run_rScript=Y</i>
--- output	
--- scripts	

3.2.2 Module Output

#	Module	Output Files										
1	SeqTrimmer	Trimmed fasta/fastq sequence files										
2	PairedSeqMerger	Merged fastq files										
3	Rarefier	Rarefied fasta/fastq sequence files										
4	QiimePreProcessor	QIIME only, Fast-A files & QIIME Mapping file										
5a.1	ClosedRefClassifier	QIIME only, multiple otu_table.biom files										
5a.2	MergeOtuTables	QIIME only, one otu_table.biom file										
5b	OpenRefClassifier	QIIME only, one otu_table.biom file										
5c	DeNovoClassifier	QIIME only, one otu_table.biom file										
6	ClassifierModule	<div>Classifier output:</div> <table><tr><td>RDP</td><td><sample_id>_reported.tsv</td></tr><tr><td>QIIME</td><td>otu_by_taxa_level/out_table_L<#>.txt (OTU Levels #1 – #7 = domain – species)</td></tr><tr><td>Kraken</td><td><sample_id>_reported.tsv</td></tr><tr><td>MetaPhlAn</td><td><sample_id>_reported.tsv</td></tr><tr><td>SLIMM</td><td><sample_id>_<taxonomy_level>_reported.tsv</td></tr></table> <div>If <i>report.numReads</i>=Y, add numReads to:<ul style="list-style-type: none"><i>metadata.file</i><i>metadata.descriptor</i></div>	RDP	<sample_id>_reported.tsv	QIIME	otu_by_taxa_level/out_table_L<#>.txt (OTU Levels #1 – #7 = domain – species)	Kraken	<sample_id>_reported.tsv	MetaPhlAn	<sample_id>_reported.tsv	SLIMM	<sample_id>_<taxonomy_level>_reported.tsv
RDP	<sample_id>_reported.tsv											
QIIME	otu_by_taxa_level/out_table_L<#>.txt (OTU Levels #1 – #7 = domain – species)											
Kraken	<sample_id>_reported.tsv											
MetaPhlAn	<sample_id>_reported.tsv											
SLIMM	<sample_id>_<taxonomy_level>_reported.tsv											
7	ParserModule	<div>For each level defined in <i>report.taxonomyLevels</i>:</div> <ul style="list-style-type: none"><taxonomy_level>_RawCount_metaMerged.txt<taxonomy_level>_LogNormal_metaMerged.txt <div>If <i>report.numHits</i> =Y, add numHits to:<ul style="list-style-type: none"><i>metadata.file</i><i>metadata.descriptor</i></div>										
8	RScriptBuilder	<div>For each level defined in <i>report.taxonomyLevels</i>:</div> <ul style="list-style-type: none">boxplots_<taxonomy_level>.pdfmeta_pValuesFor_<taxonomy_level>.txt										

4. QIIME Modules

- Detailed QIIME script descriptions available online: qiime.org/scripts
- Module 4: Qiime Preprocessor prepares input for any OTU picking method
- Property *qiime.pickOtus* determines the OTU picking method (5a/5b/5c)
- Module 5a.1: Closed Reference Classifier batches input for parallel processing
- Module 5a.2: Merge OTU Tables collates results from Module 5a.1
- Module 5b: Open Reference Classifier picks closed reference OTUs and attempts to classify the remaining sequences via the de novo method (default UCLUST)
- Module 5c: De Novo Classifier picks OTUs via clustering algorithm (default UCLUST)
- Module 6: Qiime Classifier builds taxonomy reports with OTUs counts by sample

4.1 Qiime Preprocessor

- This module prepares input files for classification by QIIME OTU picking script
- Decompress gzipped files and convert fastq to fasta format, if needed
- Create QIIME mapping by adding and/or reordering *metadata.file* columns
- If adding fields to *metadata.file*, assign categorical data type in *metadata.descriptor*

Module #	4	
Java Class	bioLockJ.module.classifier.r16s.qiime.QiimePreprocessor.java	
QIIME Scripts	print_qiime_config.py	Print version information to qsub output file
	validate_mapping_file.py	Build and/or verify QIIME mapping file
Control Properties	<i>control.runClassifier</i>	Required value = Y
	<i>qiime.preprocessInput</i>	Required value = Y
	<i>qiime.formatMetadata</i>	If value = Y, build qiimeMapping.txt
Input	Fasta or Fastq forward reads, which may be gzipped (paired reads must be merged via PairedSeqMerger module)	

Temp Directory	File	Description
	*.fasta or *.fastq	Decompressed gzipped sequence files
	qiimeMapping.txt	If <i>qiime.formatMetadata=Y</i> , add QIIME mapping fields in <i>metadata.file</i> if needed
	orderedMapping.txt	If <i>qiime.formatMetadata=Y</i> , reorder required columns in <i>metadata.file</i> if needed
Output Directory	File	Description
	.fasta	Primary output (Fast-A files) QiimeClassifier module pick_.py input
	mapping/*_corrected.txt	validate_mapping_file.py output
	<i>metadata.descriptor</i>	If <i>qiime.formatMetadata=Y</i> , add new fields to descriptor file if needed

4.2 Closed Reference Classifier

- This module batches input for parallel processing based on *script.batchSize*
- Each batch contains a subset of fasta files and a QIIME mapping (batchMapping.txt)
- Pick closed reference OTUs from a reference database for each batch
- Each batch outputs classifies reads by Sample ID into its own otu_table.biom file
- The MergeOtuTables module must run next to combine otu_table.biom files

Module #	5a.1	
Java Class	bioLockJ.module.classifier.r16s.qiime.ClosedRefClassifier.java	
QIIME Scripts	add_qiime_labels.py	Build combined_seqs.fna
	pick_closed_reference_otus.py	Build batch otu_table.biom files
Control Properties	<i>control.runClassifier</i>	Required value = Y
	<i>qiime.pickOtu</i>	Required value = Y
	<i>qiime.pickOtuScript</i>	Required value = pick_closed_reference_otus.py

Input	Fast-A sequence files	
Temp Directory	N/A	
Output Directory	A <i>batch_#</i> directory for every <i>script.batchSize</i> # fasta files contains:	
	File	Description
	otu_table.biom	pick_closed_reference_otus.py output MergeOtuTables module input
	97_otus.tree	pick_closed_reference_otus.py output
	batchMapping.txt	QIIME mapping for fasta/* files
	combined_seqs.fna	add_qiime_labels.py output
	fasta/*.fasta	<i>script.batchSize</i> # fasta files
	log_*.txt	pick_closed_reference_otus.py log file
	uclust_ref_picked_otus/*	pick_closed_reference_otus.py output

4.3 Merge OTU Tables

- This module combines otu_table.biom files output by Module 5a.1 ClosedRefClassifier

Module #	5a.2	
Java Class	bioLockJ.module.classifier.r16s.qiime.MergeOtuTables.java	
QIIME Scripts	merge_otu_tables.py	Merge ClosedRefClassifier otu_table.biom files
Control Properties	<i>control.runClassifier</i>	Required value = Y
	<i>qiime.mergeOtuTables</i>	Required value = Y
	<i>qiime.pickOtuScript</i>	Required value = pick_closed_reference_otus.py
Input	ClosedRefClassifier/output/batch_*/otu_table.biom files	
Temp Directory	N/A	

Output Directory	File	Description
	otu_table.biom	merge_otu_tables.py output QiimeClassifier module input

4.4 Open Reference Classifier

- QIIME picks closed reference OTUs from a reference database (see Section 5.2)
- Unclassified reads are clustered with QIIME de novo method (see Section 5.5)
- If *qiime.removeChimeras=Y*, vsearch is used to find chimeric sequences
- QIIME script filter_otus_from_otu_table.py is used to remove chimeric sequences

Module #	5b	
Java Class	bioLockJ.module.classifier.r16s.qiime.OpenRefClassifier.java	
QIIME Scripts	add_qiime_labels.py	Build combined_seqs.fna
	pick_open_reference_otus.py	Build otu_table_*.biom files
	filter_otus_from_otu_table.py	Filter chimeras from otu_table_*.biom Build primary output: otu_table.biom
Control Properties	<i>control.runClassifier</i>	Required value = Y
	<i>qiime.pickOtus</i>	Required value = Y
	<i>qiime.pickOtuScript</i>	Required value = pick_open_reference_otus.py
	<i>qiime.removeChimeras</i>	If value = Y, remove chimeras with vsearch
Input	Fast-A sequence files	
Temp Directory	N/A	

Output Directory	File	Description
	otu_table.biom	Primary output QiimeClassifier module input
	chimeras.fasta	vsearch output filter_otus_from_otu_table.py input
	combined_seqs.fna	add_qiime_labels.py output
	final_otu_map*.txt	pick_open_reference_otus.py output
	index.html	pick_open_reference_otus.py output
	log_*.txt	pick_open_reference_otus.py log file
	new_refseqs.fna	pick_open_reference_otus.py output
	nochimeras.fasta	vsearch output
	otu_table_mc*.biom	pick_open_reference_otus.py output filter_otus_from_otu_table.py input
	pynast_aligned_seqs/*	pick_open_reference_otus.py output
	rep_set.fna	pick_open_reference_otus.py output vsearch input
	rep_set.tre	pick_open_reference_otus.py output
	step*_otus/*	pick_open_reference_otus.py output
	uclust_assigned_taxonomy/*	pick_open_reference_otus.py output

4.4.1 QIIME Open Reference [index.html]

- The Open Reference index.html report provides a summary of output files



Run summary data	
Run summary data	log_20170708144403.txt
Taxonomy assignments	
OTU taxonomic assignments	rep_set_tax_assignments.txt
OTU tables	
OTU table excluding OTUs with fewer than 2 sequences	otu_table_mc2.biom
OTU table excluding OTUs with fewer than 2 sequences and including OTU taxonomy assignments	otu_table_mc2_w_tax.biom
OTU table excluding OTUs with fewer than 2 sequences and sequences that fail to align with PyNAST and including OTU taxonomy assignments	otu_table_mc2_w_tax_no_pynast_failures.biom
Trees	
OTU phylogenetic tree	rep_set.tre
OTU maps	
Final map of OTU identifier to sequence identifiers excluding OTUs with fewer than 2 sequences	final_otu_map_mc2.txt
Sequences	
OTU representative sequences	rep_set.fna
New reference sequences (i.e., OTU representative sequences plus input reference sequences)	new_refseqs.fna

4.5 De Novo Classifier

- QIIME uses the de novo OTU picking script to cluster reads and assign taxonomy
- If *qiime.removeChimeras=Y*, vsearch is used to find chimeric sequences
- QIIME script filter_otus_from_otu_table.py is used to remove chimeric sequences

Module #	5c		
Java Class	bioLockJ.module.classifier.r16s.qiime.DeNovoClassifier.java		
QIIME Scripts	add_qiime_labels.py		Build combined_seqs.fna
	pick_open_reference_otus.py		Build otu_table_*.biom files
	filter_otus_from_otu_table.py		Filter chimeras from otu_table_*.biom Build primary output: otu_table.biom
Control Properties	control.runClassifier		Required value = Y
	qiime.pickOtu		Required value = Y
	qiime.pickOtuScript		Required value = pick_de_novo_otus.py
	qiime.removeChimeras		If value = Y, remove chimeras with vsearch
Input	Fast-A sequence files		
Temp Directory	N/A		
Output Directory	Fast-A files must all be processed together to find de novo OTUs		
	File		Description
	otu_table.biom		Primary output QiiimeClassifier module input
	chimeras.fasta		vsearch output filter_otus_from_otu_table.py input
	combined_seqs.fna		add_qiime_labels.py output
	log_*.txt		pick_de_novo_otus.py log file
	nochimeras.fasta		vsearch output
	pynast_aligned_seqs/*		pick_de_novo_otus.py output
	rep_set/*		pick_de_novo_otus.py output vsearch input
	rep_set.tre		pick_de_novo_otus.py output
	uclust_assigned_taxonomy/*		pick_de_novo_otus.py output
	uclust_picked_otus/*		pick_de_novo_otus.py output

4.6 QIIME Classifier

- QIIME Classifier builds taxonomy level reports by counting OTUs for each sample
- Alpha metrics can be included in R script by adding to *report.attributes*, if configured

Module #	6	
Java Class	bioLockJ.module.classifier.r16s.QiimeClassifier.java	
QIIME Scripts	summarize_taxa.py	Create otu_by_taxa_level/*
	alpha_diversity.py	Create alphaDiversity.txt
	add_alpha_to_mapping_file.py	Add alpha metrics to QIIME mapping
Control Properties	<i>control.runClassifier</i>	Required value = Y
	<i>project.classifierType</i>	Required value = QIIME
	<i>qiime.alphaDiversityMetrics</i>	If <i>qiime.alphaDiversityMetrics</i> configured add alpha metrics to qiimeMapping.txt and <i>metadata.descriptor</i>
Input	otu_table.biom	
Temp Directory	N/A	
Output Directory	File	Description
	otu_by_taxa_level/*	summarize_taxa.py output
	alphaDiversity.txt	alpha_diversity.py output
	otuSummary.txt	“biom summarize-table” output
	qiimeMapping.txt	add_alpha_to_mapping_file.py output
	<i>metadata.descriptor</i>	Add <i>qiime.alphaDiversityMetrics</i> fields to descriptor file, if any

5. Extending the Pipeline

- There are several options to extend the pipeline to meet your project needs:
 - Add complex filters or multi-variate statistical models to the R Script
 - Run datasets through multiple classifiers to compare accuracy and performance
 - Add new Java modules to support your favorite classifier

5.1 Enhancing the R Script

- R Script Module identifies significant OTUs for individual metadata *report.attributes*
- Generate new reports by updating report.r with complex filters and new models

5.2 Comparing Classifiers

- Update *project.classifierType* & add the [classifier specific properties](#) to switch classifiers
 - Compare significant OTUs output by RScriptBuilder
 - Compare raw count and relative abundance tables output by ParserModule
 - Compare overall performance with the BioLockJ runtime summary

5.3 Adding New Classifiers

- Add support for additional classifiers by extending 3 abstract Java classes:
 1. ClassifierModule.java
 2. ParserModule.java
 3. OtuNode.java

5.3.1 Extending ClassifierModule.java

- Validate new classifier properties in checkDependencies()
- Implement abstract methods that build bash scripts for new classifier:
 - protected abstract** List<List<String>> buildScript(**final** List<File> files)
 - protected abstract** List<List<String>> buildScriptForPairedReads(**final** List<File> files)
- Each inner List<String> holds a group of statements used to classify one sequence file

5.3.2 Extending ParserModule.java

- Implement abstract method to build BioLockJ OTU Nodes:
 - protected void** createOtuNodes()
- Parse classifier output and instantiate new OtuNode subclass for each line
- Call addOtuNode(id, node) to map Sample ID to OTU count

5.3.3 Extending OtuNode.java

- Constructor accepts a String value representing one line of classifier output
- Parse line for OTU count

Appendix A: Contact Information

BioLockJ Administrator:	biolockj@gmail.com
Project URL:	github.com/mikesioda/BioLockJ
Principal Investigator:	Dr. Anthony Fodor <afodor@uncc.edu>
Lead Developer:	Michael Sioda <msioda@uncc.edu>