

**INVESTIGATING TRANSCRIPTOMIC RESPONSES TO BIOFUEL
STRESS IN *CLOSTRIDIUM ACETOBUTYLICUM*:
TRANSCRIPTOME ASSEMBLY AND GENOME ANNOTATION OF A
MODEL FERMENTATIVE BACTERIUM.**

by

Matthew T. Ralston

A thesis submitted to the Faculty of the University of Delaware in partial fulfillment of the requirements for the degree of Master of Science in Bioinformatics and Computational Biology

Summer 2014

© 2014 Matthew T. Ralston
All Rights Reserved

INVESTIGATING TRANSCRIPTOMIC RESPONSES TO BIOFUEL
STRESS IN *CLOSTRIDIUM ACETOBUTYLICUM*:
TRANSCRIPTOME ASSEMBLY AND GENOME ANNOTATION OF A
MODEL FERMENTATIVE BACTERIUM.

by

Matthew T. Ralston

Approved: _____
Eleftherios T. Papoutsakis, Ph.D.
Professor in charge of thesis on behalf of the Advisory Committee

Approved: _____
Errol Lloyd, Ph.D.
Chair of the Department of Computer Science

Approved: _____
Babatunde Ogunnaike, Ph.D.
Dean of the College of Engineering

Approved: _____
James G. Richards, Ph.D.
Vice Provost for Graduate and Professional Education

ACKNOWLEDGMENTS

I write this paper with unending thanks for my family, friends, the love, support and encouragement that they give, and the lessons they have taught me. With special thanks for my mother Donna, father Thomas, and sister Allison. With special thanks to my mother; her character and duty towards others is inspirational for this work's focus on sustainability and climate change. With special thanks to my Father, for inspiration through his work ethic, leadership, and open mind. With special thanks to my sister for her support of my education, curiosity, and maturity; I would not be where I am without her encouragement and acceptance. I write this with gratitude to my family for the life they have given me, the sacrifices they have made on my behalf, and most importantly their love. With special thanks for my grandfather George inspiring my pursuit of science and chemistry. With special thanks for my grandmother Winnie for the inspiration of her love and optimism which bring me through each challenge. With special thanks for my uncle Kevin for his curiosity, friendship, optimism, and support. With special thanks for my brother-in-law Matt for his friendship and encouragement. With special thanks for my nieces Violet and Ruby for their love, for the lessons that they teach me, and their infectious energy and optimism. Many thanks for Madeline for her love, open ear, curiosity, support, and encouragement throughout the years. Many thanks to my best friend Andrew for his friendship, curiosity, support, and encouragement. With many thanks to the rest of my wonderful and supportive family and friends for their unconditional love and support.

Many thanks to Karol Miaskiewicz for his consistent and wonderful friendship throughout this project. Many thanks to all of the members of the Bioinformatics program. Particularly, I'd like to thank Erin Crowgey for her mentorship and support of my efforts to learn NGS bioinformatic analyses and Ryan Moore for tossing ideas

around together. I'd like to thank Shawn Polson for his open ear and perspective. I'd also like to thank Dr. Wu for her support and encouragement throughout the program and the many members of the Wu group for their support as well. Many thanks to Bruce, Olga, and Summer for their work in the Sequencing and Genotyping Center in support of this project.

Many thanks for my mentors, Drs. Keerthi Venkataramanan and Terry Papoutsakis, for their mentorship, support, critique, and understanding throughout this project. The success that we've seen in this effort I owe to Keerthi's guidance and training. Many thanks for each of the many members of the Papoutsakis lab for their friendship, support, and critique. I am very grateful for the unity of this group and for how much I enjoy coming to the lab each day. For the hard work, sleepless nights, early autoclave cycles, shared spaces, and plenty of reasons to celebrate I am grateful for each member.

Ad maiorem dei gloriam.

TABLE OF CONTENTS

LIST OF TABLES	vi
LIST OF FIGURES	vii
ABSTRACT	viii

Chapter

1 METHODS	1
1.1 Culture	1
1.2 RNA preparation	1
1.3 RNA enrichment, RNA-seq library preparation, and Sequencing . . .	2
1.4 Sequencing Statistics, Alignments	4
1.5 Transcriptome Assembly	5
1.6 Coverage and Differential Expression	5
1.7 Molecular Methods	6
2 RESULTS	7
3 DISCUSSION	8
4 CONCLUSIONS	9

LIST OF TABLES

LIST OF FIGURES

ABSTRACT

Chapter 1

METHODS

1.1 Culture

Wild type *Clostridium acetobutylicum* ATCC 824 was cultured anaerobically in 4L New Brunswick Scientific BioFlo 310 bioreactors at 37 °C, pH 5.0, 200 mL min⁻¹ N₂ and 200rpm agitation in a defined *Clostridia* growth medium, as described previously¹. When the cultures were grown to A₆₀₀=1, the N₂ flow rate was decreased to 50 mL min⁻¹ and cultures were either stressed to a final concentration of 60 mM *n*-butanol, 40 mM potassium butyrate, or left unstressed. 15 mL samples were acquired at 15, 75, 150, and 270 minutes after treatment and OD synchronization. Samples were centrifuged at 8,000rpm, 4 °C for 20 minutes. After discarding the supernatant, cell pellets were then immediately frozen at -85 °C.

1.2 RNA preparation

RNA was extracted by first washing the cell pellets in 1mL of RNase-free SET buffer (25% sucrose, 50 mM EDTA [pH 8.0], 50 mM Tris-HCl [pH 8.0]) before resuspending cells in a 220 mL solution of RNase-free SET buffer containing 4.55 U mL⁻¹ proteinase K and 20 mg mL⁻¹ lysozyme and incubating for 6 minutes. Resuspended cells were vortexed with 40mg of RNase-free glass beads ($\leq 106 \mu\text{m}$) at maximum speed and room temperature for 4 minutes. Each sample was mixed immediately with 1 mL of ice-cold QIAzol (Qiagen, Valencia, CA, USA) and then 200 μL of ice-cold chloroform, mixing well with each addition. After a 3 minute room temperature incubation, samples were centrifuged at 11,000rpm and 4 °C for 15 minutes. The aqueous phase was then mixed with 1.3 mL of ice-cold ethanol before transferring to a miRNeasy Mini

spin-column (Qiagen, Valencia, CA, USA) and centrifuging at 11,000rpm and 4 °C for 15 seconds.

Next, 700 μ L of RWT buffer was added to the column, before centrifuging at 11,000rpm and 4 °C for 15 seconds, discarding the collection tube and transferring the column to a fresh collection tube. The column was washed twice with 500 μ L of RPE buffer before centrifuging at 11,000rpm and 4 degreeCelsius for 15 seconds each. The membrane was then dried with an additional centrifugation step at 11,000rpm and 4 degreeCelsius for 1 minute. The RNA was eluted twice by incubating with 50 μ L of nuclease-free water for 1 minute and eluting for 1 minute at 11,000rpm and 4 degreeCelsius.

After quantification on a Nanodrop ND-1000, samples were then precipitated in 0.3M sodium acetate and 75% ethanol overnight, centrifuged at 14,000 rpm for 30 minutes, washed twice with 400 μ L ice-cold 70% ethanol, and rehydrated in 50 μ L RNase-free water. Next, samples were treated with the Turbo DNA-free kit (Ambion, Austin, TX, USA). 5 μ L of 10X Turbo DNase buffer and 1 μ L of Turbo DNase ($2\text{U}\mu\text{L}^{-1}$) were added to each sample before incubating at 37 degreeCelsius for 30 minutes. Next, 5 μ L of DNase inactivation reagent were added to each sample, mixing occasionally for 5 minutes. The samples were then centrifuged at 10,000rpm and 4 degreeCelsius for 90 seconds, precipitating the DNase. The samples were moved to fresh 1.5 μ L tubes.

Samples were then precipitated, washed twice more with 70% ethanol, and resuspended in 20 μ L of nuclease-free water, requantified, and aliquoted for quality analysis with the BioAnalyzer platform (Agilent, Wilmington, DE, USA), and 10 μ g aliquots in 10 μ L samples were stored at -85°C .

1.3 RNA enrichment, RNA-seq library preparation, and Sequencing

Ribosomal RNA was removed with the MicrobExpress kit (Ambion, Austin, TX, USA) according to their protocol. Briefly, beads were prepared by taking 50 μ L for each sample, washing with an equal volume (50 μ L) of water capturing for 5 minutes on a MagnaSphere (Promega, Madison, WI, USA) magnetic stand and aspirating.

Subsequently, the beads were resuspended in an equal volume (50 μ L each) of binding buffer and capturing as above. The beads were then resuspended in an equal volume (50 μ L each) of binding buffer and warmed to 37 °C. Next, 200 μ L of binding buffer was added to each 10 μ g RNA aliquot with 4 μ L of capture oligo mix. The mixture was warmed to 70 °C for 10 minutes, then cooled to 37 °C for 15 minutes. Next, the rRNA was captured by mixing 50 μ L of beads with each sample, incubating for 15 minutes at 37 °C, and capturing as above. The enriched RNA was transferred to a fresh 1.5 mL tube. The beads were then washed with 100 μ L of pre-warmed (37 °C) wash solution, incubating on the magnetic stand for 5 minutes, and adding the wash solution to the enriched RNA. The samples were then ethanol precipitated at 20 °C overnight with 35 μ L of 3 M Sodium Acetate, 5 mg mL⁻¹ Glycogen, and 1175 μ L of chilled 100% ethanol. The samples were washed twice with 70% ethanol and resuspended in 25 μ L. The samples were enriched further by repeating the MicroExpress treatment. Small 10-100 ng aliquots were analyzed at each step with the BioAnalyzer to monitor enrichment.

Selected samples were enriched further with Terminator 5'-phosphate dependent exonuclease kit (Epicentre, Madison, WI, USA). Terminator Exonuclease 1 μ L (1 U μ L⁻¹) was added with 2 μ L 10X Buffer A to each RNA sample. The reaction was run in a thermocycler for 60 minutes at 30 °C. The reaction was terminated with the addition of 1 μ L of 100 mM EDTA and 100 mM Tris HCl at pH 8.0. The samples were then purified by ethanol precipitation (0.3 M Sodium Acetate and 75% ethanol) with two 70% ethanol washes, as above. Enriched RNA was quantified as above and assessed for quality with the BioAnalyzer platform (Agilent, Wilmington, DE, USA). High quality samples were used to prepare RNA-seq libraries with the ScriptSeq v2 library preparation kit and indexed PCR primers (Epicentre, Madison, WI, USA). Briefly, 1 μ L of fragmentation solution and 2 μ L of cDNA synthesis primer was added to 50 ng of RNA and the solution was fragmented for 5 minutes at 85 °C in a thermocycler. To each reaction, 0.5 mM of Dithiothreitol, 3 μ L of cDNA synthesis premix, 0.5 μ L StarScript Reverse Transcriptase. is added to each sample and run with the following cycle: 5

minutes at 25 °C, 20 minutes at 42 °C. After cooling each reaction to 37 °C, 1 µL of finishing solution was added, incubating for 10 minutes. The RNA is degraded by fragmenting further for 3 minutes at 95 °C, cooling to 25 °C. The first strand cDNA is di-tagged by adding 7.5 µL of terminal tagging premix and 0.5 µL of DNA polymerase. The terminal tagging reaction is run at 25 °C for 15 minutes and 95 °C for 3 minutes. The di-tagged cDNA is then purified with the AMPure XP bead system (Beckmann Coulter, Brea, CA, USA). First, the library is mixed with 45 µL of homogenous bead mixture. After thorough mixing, each solution is transferred to a 1.5 mL tube and the library is captured with the magnetic stand and the supernatant aspirated. Each library is then washed twice with 200 µL of 80% ethanol. After resuspending in 24.5 µL of nuclease-free water, the beads are captured and each library is transferred to a new 200 µL microfuge tube. Adapters are added to the di-tagged cDNA during PCR by adding 25 µL FailSafe Premix E, 1 µL forward primer, 1 µL of ScriptSeq v2 indexed reverse PCR primer, 0.5 µL of FailSafe Polymerase. The PCR conditions are as follows: cycles of 30 seconds of 95 °C, 30 seconds of 55 °C, and 3 minutes of 68 °C. After 12 cycles, the reaction terminates with a 7 minute incubation at 68 °C before purifying the library with the AMPure system, as above. Libraries were multiplexed and sequenced for 76 cycles over two lanes of an Illumina HiSeq 2500 at the University of Delaware Sequencing and Genotyping Center (Newark, DE, USA).

1.4 Sequencing Statistics, Alignments

Paired-end sequencing resulted in 1000 pairs of 76 bp reads which are deposited in the Sequence Read Archive (SRA). Summary statistics for the libraries are shown in table/appendix (Table 1). The basic bioinformatic processing pipeline is described on [Github](#). In brief, the fastq headers are briefly pre-processed for downstream applications. Then, adaptors are removed from the reads with Trimmomatic²super. Base quality is adjusted by trimming to the minimum base quality of 20. Before aligning to the *Clostridium acetobutylicum* ATCC 824 genome, the data is subjected to *in silico* ribosomal RNA removal. The reads are aligned to the rRNA sequences with

Bowtie 2.1.0³. The unmapped reads are then aligned to the genome and megaplasmid sequences (NC'003030.1 and NC'001988.2). The alignment files were then cleaned, sorted, indexed, and validated with SAMtools⁴ and Picard⁵.

1.5 Transcriptome Assembly

Reference and *de novo* assembly was done with Trinity⁶. Fastq files were modified by appending the second column of the fastq Casava 1.8+ header to the first before processing and alignment. Next, the resulting alignment files were merged and sorted before appending the pair information (/1 or /2) according to the Trinity documentation. To assess the assemblies, I have been contributing to a transcriptome assembly assessment software project: [Transrate](#). This software assesses transcriptome assemblies by calculating coverage statistics and assessing the assemblies agreement with the reference proteome. I have made several additions to this software. Specifically, unpaired reads and strand specific alignment were integrated into the coverage/alignment statistics. Additionally, singleton reads were produced from the alignment process and were then further assessed for possible sources of contamination. Finally, the assembly itself was aligned to the reference genome, assuring the validity of the assembly and the identity of the assembled transcripts.

1.6 Coverage and Differential Expression

Coverage vectors for each strand were calculated with BEDtools⁷. Coverage vectors for each transcript were then acquired with a custom Ruby script. Summarization and visualization of these data was performed in R⁸. Next, transcriptional start sites were identified with TSSi⁹ after extracting (200? 300?) bp windows surrounding the beginning of the assembled transcript with a custom Ruby script. Coverage vectors from each transcript were visualized with these results in R. Read counts per transcript were quantified with HTSeq¹⁰. Differential expression analyses followed the conservative approach of DESeq2^{11, 12}. Calculations and visualizations were done in R with various packages for heatmaps, principal components analysis, clustering, and

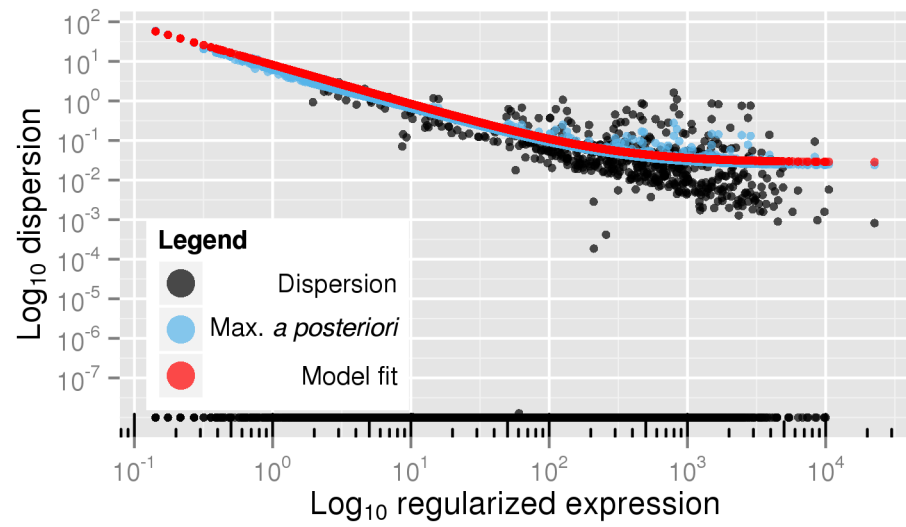
advanced graphics¹³. Data were also processed manually for visualization in Circos graphs¹⁴.

1.7 Molecular Methods

Fold changes were confirmed for randomly selected genes by qRT-PCR. Transcriptional start sites were verified by 5'RACE for randomly selected genes. Small RNAs were confirmed by RT-PCR and Northern blot.

Chapter 2

RESULTS



Chapter 3
DISCUSSION

Chapter 4
CONCLUSIONS

Bibliography

- [1] Keerthi P Venkataramanan, Shawn W Jones, Kevin P McCormick, Sridhara G Kunjeti, Matthew T Ralston, Blake C Meyers, and Eleftherios T Papoutsakis. The clostridium small rnome that responds to stress: the paradigm and importance of toxic metabolite stress in c. acetobutylicum. *BMC genomics*, 14(1):849, 2013.
- [2] Anthony M Bolger, Marc Lohse, and Bjoern Usadel. Trimmomatic: a flexible trimmer for illumina sequence data. *Bioinformatics*, page 170, 2014.
- [3] Ben Langmead and Steven L Salzberg. Fast gapped-read alignment with bowtie 2. *Nature methods*, 9(4):357–359, 2012.
- [4] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, and Richard Durbin. The sequence alignment/map format and samtools. *Bioinformatics*, 25(16):2078–2079, 2009.
- [5] Alec W. Picard, 2009.
- [6] Manfred G Grabherr, Brian J Haas, Moran Yassour, Joshua Z Levin, Dawn A Thompson, Ido Amit, Xian Adiconis, Lin Fan, Raktima Raychowdhury, Qiandong Zeng, et al. Full-length transcriptome assembly from rna-seq data without a reference genome. *Nature biotechnology*, 29(7):644–652, 2011.
- [7] Aaron R Quinlan and Ira M Hall. Bedtools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, 26(6):841–842, 2010.
- [8] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2014.

- [9] Clemens Kreutz, JS Gehring, D Lang, Ralf Reski, Jens Timmer, and Stefan A Rensing. Tssian r package for transcription start site identification from 5 mrna tag data. *Bioinformatics*, 28(12):1641–1642, 2012.
- [10] Simon Anders. Htseq: Analysing high-throughput sequencing data with python. URL <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>, 2010.
- [11] Michael I Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *bioRxiv*, 2014.
- [12] Simon Anders and Wolfgang Huber. Differential expression analysis for sequence count data. *Genome biol*, 11(10):R106, 2010.
- [13] Hadley Wickham. *ggplot2: elegant graphics for data analysis*. Springer, 2009.
- [14] et al. Krzywinski, Martin. Circos: an information aesthetic for comparative genomics. *Genome research*, 19(9):1639–1645, 2009.
- [15] Keith V Alsaker, Thomas R Spitzer, and Eleftherios T Papoutsakis. Transcriptional analysis of spo0a overexpression in clostridium acetobutylicum and its effect on the cell’s response to butanol stress. *Journal of bacteriology*, 186(7):1959–1971, 2004.
- [16] S Andrews. Fastqc: A quality control tool for high throughput sequence data. *Reference Source*, 2010.
- [17] Dominik Antoni, Vladimir V Zverlov, and Wolfgang H Schwarz. Biofuels from microbes. *Applied microbiology and biotechnology*, 77(1):23–35, 2007.
- [18] Shota Atsumi, Anthony F Cann, Michael R Connor, Claire R Shen, Kevin M Smith, Mark P Brynildsen, Katherine JY Chou, Taizo Hanai, and James C Liao. Metabolic engineering of *escherichia coli* for 1-butanol production. *Metabolic engineering*, 10(6):305–311, 2008.

- [19] Timothy L Bailey, Mikael Boden, Fabian A Buske, Martin Frith, Charles E Grant, Luca Clementi, Jingyuan Ren, Wilfred W Li, and William S Noble. Meme suite: tools for motif discovery and searching. *Nucleic acids research*, 37(suppl 2):W202–W208, 2009.
- [20] Jacob R Borden and Eleftherios Terry Papoutsakis. Dynamics of genomic-library enrichment and identification of solvent tolerance genes for *clostridium acetobutylicum*. *Applied and environmental microbiology*, 73(9):3061–3068, 2007.
- [21] Yili Chen, Dinesh C Indurthi, Shawn W Jones, and Eleftherios T Papoutsakis. Small rnas in the genus *clostridium*. *MBio*, 2(1):e00340–10, 2011.
- [22] Patrik Dhaeseleer, Shoudan Liang, and Roland Somogyi. Genetic network inference: from co-expression clustering to reverse engineering. *Bioinformatics*, 16(8):707–726, 2000.
- [23] A Gordon and GJ Hannon. Fastx-toolkit. *FASTQ short-reads pre-processing tools*, 2010.
- [24] Michael Hecker, Sandro Lambeck, Susanne Toepfer, Eugene Van Someren, and Reinhard Guthke. Gene regulatory network inference: data integration in dynamic modelsa review. *Biosystems*, 96(1):86–103, 2009.
- [25] Michael Hecker, Wolfgang Schumann, and Uwe Vlker. Heatshock and general stress response in *bacillus subtilis*. *Molecular microbiology*, 19(3):417–428, 1996.
- [26] Eric CH Ho, Michael E Donaldson, and Barry J Saville. *Detection of antisense RNA transcripts by strand-specific RT-PCR*, pages 125–138. Springer, 2010.
- [27] Jay D Keasling, Abraham Mendoza, and Phil S Baran. Synthesis: A constructive debate. *Nature*, 492(7428):188–189, 2012.
- [28] Donghyuk Kim, Jay Sung-Joong Hong, Yu Qiu, Harish Nagarajan, Joo-Hyun Seo, Byung-Kwan Cho, Shih-Feng Tsai, and Bernhard Palsson. Comparative

- analysis of regulatory elements between *escherichia coli* and *klebsiella pneumoniae* by genome-wide transcription start site profiling. *PLoS genetics*, 8(8):e1002867, 2012.
- [29] Ethan I Lan and James C Liao. Metabolic engineering of cyanobacteria for 1-butanol production from carbon dioxide. *Metabolic engineering*, 13(4):353–363, 2011.
- [30] Joshua Z Levin, Moran Yassour, Xian Adiconis, Chad Nusbaum, Dawn Anne Thompson, Nir Friedman, Andreas Gnirke, and Aviv Regev. Comprehensive comparative analysis of strand-specific rna sequencing methods. *Nature methods*, 7(9):709–715, 2010.
- [31] Ronny Lorenz, Stephan HF Bernhart, Christian Hoener Zu Siederdisen, Hakim Tafer, Christoph Flamm, Peter F Stadler, and Ivo L Hofacker. Viennarna package 2.0. *Algorithms for Molecular Biology*, 6(1):26, 2011.
- [32] Marloes H Maathuis, Diego Colombo, Markus Kalisch, and Peter Bhlmann. Predicting causal effects in large-scale systems from observational data. *Nature Methods*, 7(4):247–248, 2010.
- [33] Florian Markowetz, Dennis Kostka, Olga G Troyanskaya, and Rainer Spang. Nested effects models for high-dimensional phenotyping screens. *Bioinformatics*, 23(13):i305–i312, 2007.
- [34] G Najoshi. Sickie-a windowed adaptive trimming tool for fastq files using quality.
- [35] Sergios A Nicolaou, Stefan M Gaida, and Eleftherios T Papoutsakis. A comparative view of metabolite and substrate stress and tolerance in microbial bioprocessing: from biofuels and chemicals, to biocatalysis and bioremediation. *Metabolic engineering*, 12(4):307–331, 2010.
- [36] Fatih Ozsolak and Patrice M Milos. Rna sequencing: advances, challenges and opportunities. *Nature reviews genetics*, 12(2):87–98, 2011.

- [37] Eleftherios T Papoutsakis. Engineering solventogenic clostridia. *Current opinion in biotechnology*, 19(5):420–429, 2008.
- [38] Christophe Pichon and Brice Felden. Small rna gene identification and mrna target predictions in bacteria. *Bioinformatics*, 24(24):2807–2813, 2008.
- [39] Yosef Prat, Menachem Fromer, Nathan Linial, and Michal Linial. Recovering key biological constituents through sparse representation of gene expression. *Bioinformatics*, 27(5):655–661, 2011.
- [40] Juan L Ramos, Estrella Duque, Mara-Trinidad Gallegos, Patricia Godoy, Mara Isabel Ramos-Gonzalez, Antonia Rojas, Wilson Tern, and Ana Segura. Mechanisms of solvent tolerance in gram-negative bacteria. *Annual Reviews in Microbiology*, 56(1):743–768, 2002.
- [41] Adam Roberts, Harold Pimentel, Cole Trapnell, and Lior Pachter. Identification of novel transcripts in annotated genomes using rna-seq. *Bioinformatics*, 27(17):2325–2329, 2011.
- [42] Tobias Sahr, Christophe Rusniok, Delphine Dervins-Ravault, Odile Sismeiro, Jean-Yves Coppee, and Carmen Buchrieser. Deep sequencing defines the transcriptional map of *l. pneumophila* and identifies growth phase-dependent regulated ncRNAs implicated in virulence. *RNA Biol*, 9(4):503–519, 2012.
- [43] Yogita Sardesai and Saroj Bhosle. Tolerance of bacteria to organic solvents. *Research in Microbiology*, 153(5):263–268, 2002.
- [44] Robert Schmieder and Robert Edwards. Quality control and preprocessing of metagenomic datasets. *Bioinformatics*, 27(6):863–864, 2011.
- [45] Olga A Soutourina, Marc Monot, Pierre Boudry, Laure Saujet, Christophe Pichon, Odile Sismeiro, Ekaterina Semenova, Konstantin Severinov, Chantal Le Bouguenec, and Jean-Yves Coppe. Genome-wide identification of regulatory

- rnas in the human pathogen clostridium difficile. *PLoS genetics*, 9(5):e1003493, 2013.
- [46] Morgane Thomas-Chollier, Matthieu Defrance, Alejandra Medina-Rivera, Olivier Sand, Carl Herrmann, Denis Thieffry, and Jacques van Helden. Rsat 2011: regulatory sequence analysis tools. *Nucleic acids research*, 39(suppl 2):W86–W91, 2011.
- [47] Christopher A Tomas, Jeffrey Beamish, and Eleftherios T Papoutsakis. Transcriptional analysis of butanol stress and tolerance in clostridium acetobutylicum. *Journal of bacteriology*, 186(7):2006–2018, 2004.
- [48] Bryan P Tracy, Shawn W Jones, Alan G Fast, Dinesh C Indurthi, and Eleftherios T Papoutsakis. Clostridia: the importance of their exceptional substrate and metabolite diversity for biofuel and biorefinery applications. *Current opinion in biotechnology*, 23(3):364–381, 2012.
- [49] Charles J Vaske, Stephen C Benz, J Zachary Sanborn, Dent Earl, Christopher Szeto, Jingchun Zhu, David Haussler, and Joshua M Stuart. Inference of patient-specific pathway activities from multi-dimensional cancer genomics data using paradigm. *Bioinformatics*, 26(12):i237–i245, 2010.
- [50] E Gerhart H Wagner. Kill the messenger: bacterial antisense rna promotes mrna decay. *Nature structural and molecular biology*, 16(8):804–806, 2009.
- [51] Qinghua Wang, Keerthi Prasad Venkataramanan, Hongzhan Huang, Eleftherios T Papoutsakis, and Cathy H Wu. Transcription factors and genetic circuits orchestrating the complex, multilayered response of clostridium acetobutylicum to butanol and butyrate stress. *BMC systems biology*, 7(1):120, 2013.
- [52] Ka Yee Yeung and Walter L. Ruzzo. Principal component analysis for clustering gene expression data. *Bioinformatics*, 17(9):763–774, 2001.

- [53] Yong Zhang, Tao Liu, Clifford A Meyer, Jrme Eeckhoutte, David S Johnson, Bradley E Bernstein, Chad Nusbaum, Richard M Myers, Myles Brown, and Wei Li. Model-based analysis of chip-seq (macs). *Genome Biol*, 9(9):R137, 2008.
- [54] Kyle A Zingaro and Eleftherios Terry Papoutsakis. Toward a semisynthetic stress response system to engineer microbial solvent tolerance. *Mbio*, 3(5):e00308–12, 2012.
- [55] Kyle A Zingaro and Eleftherios Terry Papoutsakis. Groesl overexpression imparts escherichia coli tolerance to, n-, and 2-butanol, 1, 2, 4-butanetriol and ethanol with complex and unpredictable patterns. *Metabolic engineering*, 15:196–205, 2013.
- [56] Helga Thorvaldsdóttir, James T Robinson, and Jill P Mesirov. Integrative genomics viewer (igv): high-performance genomics data visualization and exploration. *Briefings in bioinformatics*, page bbs017, 2012.
- [57] Eric P Nawrocki and Sean R Eddy. Infernal 1.1: 100-fold faster rna homology searches. *Bioinformatics*, 29(22):2933–2935, 2013.
- [58] Wei Zheng, Lisa M Chung, and Hongyu Zhao. Bias detection and correction in rna-sequencing data. *Bmc Bioinformatics*, 12(1):290, 2011.
- [59] Shawn ONeil and Scott Emrich. Assessing de novo transcriptome assembly metrics for consistency and utility. *BMC genomics*, 14(1):465, 2013.