1  **Title:** Shannon Entropy as a metric for conditional gene expression in *Neurospora crassa*

2

3  **Authors:** Abigail J. Ameri[1] and Zachary A. Lewis[1,2]

4

5  **Institutional Affiliations:**

6  [1] Department of Microbiology, University of Georgia, Athens, GA, 30602

7  [2] Corresponding Author: Zachary A. Lewis, University of Georgia, Athens, GA zlewis@uga.edu

8

9  ORCID

10  Abigail Ameri: https://orcid.org/0000-0001-5319-2630

11  Zachary Lewis: https://orcid.org/0000-0002-1735-8266

12

13

14

15

16

17

18

19

20

21

22

23

24

25      **Short Title:** A metric for variable gene expression

26      **Key words or phrases:** Shannon entropy, R, conditional gene expression, *Neurospora*

27      **Corresponding Author:** Zachary A. Lewis, University of Georgia, Athens, GA zlewis@uga.edu

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47  **Abstract**

48      *Neurospora crassa* has been an important model organism for molecular biology and

49  genetics for over 60 years. *N. crassa* has a complex life cycle, with over 28 distinct cell types and

50  is capable of transcriptional responses to many environmental conditions including nutrient

51  availability, temperature, and light. To quantify variation in *N. crassa* gene expression, we

52  analyzed public expression data from 97 conditions and calculated the Shannon Entropy value

53  for *Neurospora's* approximately 11,000 genes. Entropy values can be used to estimate the

54  variability in expression for a single gene over a range of conditions and be used to classify

55  individual genes as constitutive or condition-specific. Shannon entropy has previously been

56  used measure the degree of tissue specificity of multicellular plant or animal genes. We use this

57  metric here to measure variable gene expression in a microbe and provide this information as a

58  resource for the *N. crassa* research community. Finally, we demonstrate the utility of this

59  approach by using entropy values to identify genes with constitutive expression across a wide

60  range of conditions and to identify genes that are activated exclusively during sexual

61  development.

62

63

64 **Introduction**

65       Across conditions, individual genes can display expression patterns that can range from

66   conditional to constitutive. When performing Quantitative Reverse Transcription PCR (qRT-PCR)

67   it is crucial to identify constitutively expressed genes for experimental normalization (HUGGETT

68   *et al.* 2005). Conversely, highly regulated, condition-specific gene promoters are often used in

69   molecular biology to drive conditional expression of a gene under investigation (e.g., an

70   essential gene) or to control expression of reporter genes in certain cell types or environmental

71   conditions (e.g., a gene encoding a fluorescent protein) (GILES *et al.* 1985; HURLEY *et al.* 2012;

72   LAMB *et al.* 2013). Moreover, identification of genes that are exclusively expressed during a

73   condition or cell-type of interest can reveal genes that are functionally important. Such genes

74   or promoters are often identified by examining gene expression across just a handful of

75   experimental conditions; however, with the increase in publicly available transcriptomics data it

76   is possible to quantify variation in gene expression across many conditions for a given organism.

77       In 1963, Claude Shannon laid the basis for information theory, and described the unit

78   known as Shannon entropy (SHANNON 1997). A simplistic definition of Shannon entropy is that it

79   describes the amount of information a variable can hold (VAJAPEYAM 2014). In our case, a

80   variable is a gene, and the information is the collection of expression values from different

81   conditions. If a gene is classified as having low entropy, then the expression values would be

82   generally consistent across different conditions or possess a low amount of information.

83   Instead, if a gene is classified as having high entropy, then the expression of this gene would be

84   highly variable across different conditions and contain a high amount of information.

85     Since entropy describes information contained in a variable there are a number of uses

86     for such a metric. Previous studies have used entropy to investigate cell and tissue specific

87     expression of genes (SCHUG *et al.* 2005), identify potential therapeutic targets (FUHRMAN *et al.*

88     2000), characterize periodicity in gene expression (LANGMEAD *et al.* 2002), identify cancerous

89     tissue samples (VAN WIERINGEN AND VAN DER VAART 2011), and make genomic comparisons

90     (MACHADO 2012). Studies using entropy have been carried out in human cell lines (NATHANIEL D.

91     HEINTZMAN *et al.* 2009), mouse (SCHUG *et al.* 2005), plants (ZHANG *et al.* 2006), yeast (TIMOTHY R.

92     LEZON *et al.* 2006), bacteria, phage, and metagenomes (AKHTER *et al.* 2013) but not yet in

93     filamentous fungi.

94     *Neurospora crassa* has a 43Mb genome encoding approximately ~11,000 genes

95     (BORKOVICH *et al.* 2004) (add Nature paper). There is a whole genome knock out collection, and

96     genetic, genomic, and epigenetic studies have been carried out with this organism for more

97     than 100 years (COLOT *et al.* 2006). Indeed, *N. crassa* has been used as a model organism for

98     epigenetics, testing fungal enzymes for biomass degradation, and circadian clock studies

99     (DUNLAP *et al.* 2007; TIAN *et al.* 2009; ARAMAYO AND SELKER 2013). As a resource for *N. crassa*

100    researchers, we generated an entropy value for most genes in the *N. crassa* genome using

101    publicly available RNA-seq data, and we validated this approach using previously published lists

102    of housekeeping or inducible genes. This resource has a number of useful applications for the

103    *N. crassa* community.

104

105

106      **Methods:**

107      **Public data collection:**

108      Entropy calculations were made for all genes in the *N. crassa* genome using public RNA-

109      seq data sets (97 conditions from a total of 173 separate sets including replicates; Table S1).

110

111      **Data Analysis:**

112      **Mapping, TPM and entropy calculations:**

113      HiSat2 (version 2.1.0) (KIM *et al.* 2019) was used to map all of the SRA accessions to the

114      NC12 genome (NCBI assembly: GCA_000182925.2) using appropriate parameters specific for

115      paired or single end sequence reads (with parameters –RNA-strandness RF or R) to produce

116      bam files which were then sorted and indexed using SAMtools (version 1.3) (LI *et al.* 2009). If

117      experiments contain replicates, the replicate bam files were merged together before obtaining

118      counts with featureCounts from Subread (version 1.6.2) (LIAO *et al.* 2014). FeatureCounts was

119      used with parameters -T exon to generate all counts at the gene level. Counts were imported

120      into R where we obtained TPM using the function calculateTPM from the R package scater

121      (MCCARTHY *et al.* 2017). This package takes in feature-level (in our case, gene-level) counts and

122      gene lengths and outputs the TPM values for each gene. TPM values were then used to

123      calculate the Shannon entropy using the R package BioQC (ZHANG *et al.* 2017). The function

124      entropySpecificity was used to calculate the entropy values for all genes in the genome.  To

125      examine specific genes sets, we converted from NCU accession numbers to gene identifiers

126      from NCBI Genome Assembly NC12 (GCA_000182925.2) and plotted the kernel density

127      estimation with rug plots.

128

129    Data Availability Statement:  All supplementary tables have been uploaded to Figshare. Table S1

130    contains SRA accession numbers, short descriptions, total reads, and mapped reads for each

131    public data set used.Calculated entropy values for all *N. crassa* genes are listed in Table S2. Lists

132    of all *N. crassa* genes used to benchmark the entropy values and generate panels in figure 2 and

133    3 are included in Table S3. Code used to generate the data in this manuscript is available

134    through github.  https://github.com/ajcourtney/entropy

135

136

137    **Results and Discussion:**

138        Shannon entropy values are useful in measuring the amount of variation in expression

139    levels across different tissues or growth conditions**.** In order to calculate Shannon entropy

140    values for all *Neurospora crassa* genes, we first compiled a list of available RNA-seq data sets

141    present in the NCBI sequence read archive (SRA) (Table S1). We selected datasets that were

142    generated with the wild type strain Oak Ridge strain background, but we used both mating

143    types. To calculate accurate entropy values, we needed to gather many observations of gene

144    expression across different conditions. We searched the SRA database (LEINONEN *et al.* 2011) for

145    *N. crassa* RNA-sequencing entries that were processed at different developmental stages or

146    grown under different conditions. In total we gathered 173 accessions, which represent 97

147    developmental or growth conditions. We then developed a pipeline to generate entropy values

148    for each gene (Figure 1A). Calculated entropy values are available in Table S2. We first mapped

149    to the NC12 *N. crassa* genome using HiSat2 (KIM *et al.* 2019) to generate bam files. The bam

150    files were then used to generate read counts for each gene in each condition using

151    featureCounts (LIAO *et al.* 2014), which assigns reads to genomic features. Once the count file

152    was created, we calculated normalized expression values using the Transcripts per Million

153    (TPM) normalization method to create a matrix of normalized expression values for all genes in

154    all conditions. We then used this expression matrix to calculate the Shannon entropy value for

155    each gene (ZHANG *et al.* 2017).  This generated entropy values for 10,300 out of 10,398 genes.

156    The remaining 98 genes had 0 read counts in all conditions, so we were unable to calculate

157    entropy. Our final entropy values range from 0.0506 to 6.599. 70% of the genes in the genome

158    possess low entropy values between 0.05 and 1 (7,180/10,300) (Figure 1B). These values

159    include the constitutively expressed genes in the genome. Entropy values above one represent

160    only 30% of the genome (3,120/10,300), corresponding to genes with more condition-specific

161    expression patterns.

162    **Validation of entropy as a measure of gene expression variation in *N. crassa*.**

163         In order to determine if entropy values are a reliable predictor of expression variability

164    in a microbe, we examined the entropy values generated here for published gene sets expected

165    to be enriched for constitutively expressed genes, or conversely, for gene sets expected to

166    contain genes with highly condition-specific expression patterns. If entropy value is a reliable

167    measure of gene expression variation across conditions, housekeeping genes should be

168    enriched for genes with low entropy values, whereas sets of conditionally-induced genes are

169    expected to be enriched for high entropy values. Two previous studies identified genes useful

170    for RT-qPCR controls in *N. crassa*. One of which published a list of 38 genes classified as

171    "housekeeping genes" based on previously generated microarray and RNA-seq datasets under

172    three different conditions (quinic acid (QA) induction, circadian gene expression profiling, and

173    light response) (HURLEY et al. 2015), and the other study identified four genes by using previous

174    transcriptomic studies and genes used in related organisms to generate candidates that were

175    validated by quantitative PCR under different conditions (CUSICK et al. 2014) (Table S1). To

176    visualize the distribution of entropy values in this set of 42 "housekeeping" genes, we plotted a

177    kernel density estimation (KDE) of entropy values (Figure 2A). The KDE is a smoothed version of

178    a histogram estimated from the underlying data. As expected, the highest density of data

179    points in the housekeeping data set is around 0.25 (low entropy) and the density falls sharply

180    around 0.75 (Figure 2A). Two genes in this set possess entropy values above 1.6 and they

181    encode an exo-beta-1,3-glucanase and a UDP-glucose dehydrogenase. We plotted a heatmap

182    depicting TPM values for each gene in each condition with genes ranked by entropy values from

183    low to high (top to bottom) (Figure 2B). Genes with higher entropy values showed significant

184    induction of gene expression under certain conditions, whereas genes with low entropy values

185    displayed consistent expression values across all conditions. In particular, the two genes with

186    high entropy values showed marked induction under certain conditions. Thus, these data

187    highlight the value of performing a comprehensive analysis of conditional gene expression

188    when selecting constitutive control genes.

189         We further validated the use of entropy as a measure for constitutive gene expression

190    by using the same approach with a published list of genes 2,624 genes involved in transcription

191    and translation (Table S1), reasoning that genes involved in these essential processes would be

192    expressed at similar levels in all 93 conditions we investigated. (BENZ et al. 2014). The

193    distribution of entropy values for transcription and translation genes resembles the distribution

194    of entropy values for housekeeping genes where the highest density is concentrated at the low

195    end of entropy values (Figure 2C). Many of the genes that possess entropy values above 1.6 are

196    either hypothetical proteins or genes associated with cellular transport or metabolism. We

197    again examined the TPM values for each gene in this set in a heatmap ranked by entropy from

198    low to high and again find mostly steady expression across conditions (Figure 2D).

199        We next asked if higher entropy values were associated with conditionally expressed

200    genes. The highest entropy values imply that a gene must only be expressed under specific

201    conditions and may only show expression in one or a few of the conditions in the entire RNA-

202    seq dataset. To confirm that higher entropy values were indeed associated with condition- or

203    tissue-specific gene expression, we created KDE plots for 513 genes induced by light (Figure 3A

204    and Table S1) and 3,259 genes that have expression changes during sexual development (Figure

205    3C and Table S1) (Wu *et al.* 2014) (Wang *et al.* 2014). In both cases, there is a shift in

206    distribution of entropy values toward higher entropy values compared to "housekeeping" or

207    "transcription and translation" genes. We examined TPM values for each gene in each condition

208    using a heatmap ranked by entropy values from low to high (top to bottom) and find that a

209    majority of genes in each gene set show variable expression across conditions, as expected

210    (Figure 3B, D). Genes that have regulation changes during perithecial (sexual) development also

211    show a shift to the right, but with retention of more low entropy genes than in the light induced

212    gene set (Figure 3C). Plotting the TPM values in an entropy ranked heatmap shows that

213    approximately half of these genes are highly expressed across many conditions and half are

214    variably expressed, corresponding to genes with lower entropy values in the density plot

215    (Figure 3D). This implies that half of these genes are not specific to sexual or vegetative cell

216    types even though they show transcriptional changes throughout development (WANG *et al.*

217    2014).

218         As a final confirmation that entropy can be used as a reliable metric to assess the

219    variation or lack of variation in gene expression levels across many conditions, we plotted the

220    expression levels of 100 genes with the highest entropy values and 100 genes with the lowest

221    entropy values. We took the $\log_2$ TPM values for all conditions (columns) and plotted them for

222    each gene in a heatmap that was clustered by gene (row) for both the top and bottom 100

223    genes. As expected, with the lowest entropy values show mostly uniform expression across all

224    conditions (Figure 4A), and genes in the high entropy group displayed highly variable and

225    condition-specific expression (Figure 4B). Together, these data demonstrate that entropy is an

226    effective tool for measuring variation in gene expression levels in a filamentous fungus.

227         The information and code generated in the course of this study could prove useful in a

228    number of ways. First, identifying genes that are induced in a certain condition and display a

229    high entropy value will help identify genes that are condition-specific. In addition, examining

230    entropy values for individual genes can be a useful approach for finding new inducible

231    promoters to use for genetic studies. Condition-specific expressed genes are good starting

232    targets to test for this purpose.  The entropy metric determined here can also be used to

233    confirm constitutive expression of genes chosen as controls for RT-PCR. In examining the

234    housekeeping genes from previously published studies it is clear that not all will function as

235    good controls under all conditions, a limitation that was discussed by Hurley and colleagues

236    (HURLEY *et al.* 2015). We combined all of their housekeeping genes together, whereas they had

237    them divided into housekeeping genes usable for different conditions in qRT-PCR (QA

238    induction, light response studies, and circadian experiments). Here we can choose genes that

239    will work across all conditions (provided the conditions were represented in the initial dataset).

240    Our approach provides a quantitative metric that can be applied to identify condition-specific

241    genes, as opposed to investigating individual datasets or using controls from previous studies

242    which may not perform as expected. In addition, this methodology is scalable; the initial

243    inclusion of more conditions will only increase the robustness of the metric produced. As more

244    data are published, more datasets can be incorporated. This approach can be used across other

245    fungi in addition to *N. crassa*, provided there are sufficient RNA-seq data publicly available.

251

252    **References**

253    Akhter, S., B. A. Bailey, P. Salamon, R. K. Aziz and R. A. Edwards, 2013 Applying Shannon's
254          information theory to bacterial and phage genomes and metagenomes. Sci Rep 3**:** 1033.
255    Aramayo, R., and E. U. Selker, 2013 Neurospora crassa, a model system for epigenetics
256          research. Cold Spring Harb Perspect Biol 5**:** a017921.
257    Benz, J. P., B. H. Chau, D. Zheng, S. Bauer, N. L. Glass *et al.*, 2014 A comparative systems analysis
258          of polysaccharide-elicited responses in Neurospora crassa reveals carbon source-specific
259          cellular adaptations. Mol Microbiol 91**:** 275-299.
260    Borkovich, K. A., L. A. Alex, O. Yarden, M. Freitag, G. E. Turner *et al.*, 2004 Lessons from the
261          genome sequence of Neurospora crassa: tracing the path from genomic blueprint to
262          multicellular organism. Microbiol Mol Biol Rev 68**:** 1-108.
263    Colot, H. V., G. Park, G. E. Turner, C. Ringelberg, C. M. Crew *et al.*, 2006 A high-throughput gene
264          knockout procedure for Neurospora reveals functions for multiple transcription factors.
265          Proc Natl Acad Sci U S A 103**:** 10352-10357.

266    Cusick, K. D., L. A. Fitzgerald, R. K. Pirlo, A. L. Cockrell, E. R. Petersen *et al.*, 2014 Selection and
267        evaluation of reference genes for expression studies with quantitative PCR in the model
268        fungus Neurospora crassa under different environmental conditions in continuous
269        culture. PLoS One 9**:** e112706.
270    Dunlap, J. C., J. J. Loros, H. V. Colot, A. Mehra, W. J. Belden *et al.*, 2007 A circadian clock in
271        Neurospora: how genes and proteins cooperate to produce a sustained, entrainable,
272        and compensated biological oscillator with a period of about a day. Cold Spring Harb
273        Symp Quant Biol 72**:** 57-68.
274    Fuhrman, S., M. J. Cunningham, X. Wen, G. Zweiger, J. J. Seilhamer *et al.*, 2000 The application
275        of shannon entropy in the identification of putative drug targets. Biosystems 55**:** 5-14.
276    Giles, N. H., M. E. Case, J. Baum, R. Geever, L. Huiet *et al.*, 1985 Gene organization and
277        regulation in the qa (quinic acid) gene cluster of Neurospora crassa. Microbiol Rev 49**:**
278        338-358.
279    Huggett, J., K. Dheda, S. Bustin and A. Zumla, 2005 Real-time RT-PCR normalisation; strategies
280        and considerations. Genes Immun 6**:** 279-284.
281    Hurley, J. H., A. Dasgupta, P. Andrews, A. M. Crowell, C. Ringelberg *et al.*, 2015 A Tool Set for
282        the Genome-Wide Analysis of Neurospora crassa by RT-PCR. G3 (Bethesda) 5**:** 2043-
283        2049.
284    Hurley, J. M., C. H. Chen, J. J. Loros and J. C. Dunlap, 2012 Light-inducible system for tunable
285        protein expression in Neurospora crassa. G3 (Bethesda) 2**:** 1207-1212.
286    Kim, D., J. M. Paggi, C. Park, C. Bennett and S. L. Salzberg, 2019 Graph-based genome alignment
287        and genotyping with HISAT2 and HISAT-genotype. Nat Biotechnol 37**:** 907-915.
288    Lamb, T. M., J. Vickery and D. Bell-Pedersen, 2013 Regulation of gene expression in Neurospora
289        crassa with a copper responsive promoter. G3 (Bethesda) 3**:** 2273-2280.
290    Langmead, C. J., C. R. McClung and B. R. Donald, 2002 A maximum entropy algorithm for
291        rhythmic analysis of genome-wide expression patterns. Proc IEEE Comput Soc Bioinform
292        Conf 1**:** 237-245.
293    Leinonen, R., H. Sugawara, M. Shumway and C. International Nucleotide Sequence Database,
294        2011 The sequence read archive. Nucleic Acids Res 39**:** D19-21.
295    Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The Sequence Alignment/Map
296        format and SAMtools. Bioinformatics 25**:** 2078-2079.
297    Liao, Y., G. K. Smyth and W. Shi, 2014 featureCounts: an efficient general purpose program for
298        assigning sequence reads to genomic features. Bioinformatics 30**:** 923-930.
299    Machado, J. A. T., 2012 Shannon Entropy Analysis of the Genome Code. Mathematical Problems
300        in Engineering 2012.
301    McCarthy, D. J., K. R. Campbell, A. T. L. Lun and Q. F. Wills, 2017 Scater: pre-processing, quality
302        control, normalization and visualization of single-cell RNA-seq data in R. Bioinformatics
303        33**:** 1179-1186.
304    Nathaniel D. Heintzman, Gary C. Hon, R. David Hawkins, Pouya Kheradpour, Alexander Stark *et*
305        *al.*, 2009 Histone modifications at human enhancers reflect global cell-type-specific gene
306        expression. Nature 459**:** 108-112.
307    Schug, J., W. P. Schuller, C. Kappen, J. M. Salbaum, M. Bucan *et al.*, 2005 Promoter features
308        related to tissue specificity as measured by Shannon entropy. Genome Biol 6**:** R33.

309  Shannon, C. E., 1997 The mathematical theory of communication. 1963. MD Comput 14**:** 306-
310       317.
311  Tian, C., W. T. Beeson, A. T. Iavarone, J. Sun, M. A. Marletta *et al.*, 2009 Systems analysis of
312       plant cell wall degradation by the model filamentous fungus Neurospora crassa. Proc
313       Natl Acad Sci U S A 106**:** 22157-22162.
314  Timothy R. Lezon, Jayanth R. Banavar, Marek Cieplak, Amos Maritan and N. V. Fedoroff, 2006
315       Using the principle of entropy maximization to infer genetic interaction networks from
316       gene expression patterns. PNAS 103**:** 19033-19038.
317  Vajapeyam, S., 2014 Understanding Shannon's Entropy metric for Information. arXiv preprint.
318  van Wieringen, W. N., and A. W. van der Vaart, 2011 Statistical analysis of the cancer cell's
319       molecular entropy using high-throughput data. Bioinformatics 27**:** 556-563.
320  Wang, Z., F. Lopez-Giraldez, N. Lehr, M. Farre, R. Common *et al.*, 2014 Global gene expression
321       and focused knockout analysis reveals genes associated with fungal fruiting body
322       development in Neurospora crassa. Eukaryot Cell 13**:** 154-169.
323  Wu, C., F. Yang, K. M. Smith, M. Peterson, R. Dekhang *et al.*, 2014 Genome-wide
324       characterization of light-regulated genes in Neurospora crassa. G3 (Bethesda) 4**:** 1731-
325       1745.
326  Zhang, J. D., K. Hatje, G. Sturm, C. Broger, M. Ebeling *et al.*, 2017 Detect tissue heterogeneity in
327       gene expression data with BioQC. BMC Genomics 18**:** 277.
328  Zhang, X., J. Yazaki, A. Sundaresan, S. Cokus, S. W. Chan *et al.*, 2006 Genome-wide high-
329       resolution mapping and functional analysis of DNA methylation in arabidopsis. Cell 126**:**
330       1189-1201.
331
332

333     Figure Legends

334     **Figure 1: Calculation of Shannon entropy for *N. crassa* genes using public RNA-seq data.**

335     A) Schematic of our computational pipeline for calculating Shannon entropy from publically

336     available datasets.

337     B) *N. crassa* genes display a broad range of entropy values. The histogram shows entropy values

338     for all genes. The y-axis is the number of genes found in each bin. The x-axis shows the binned

339     entropy values.

340

341     **Figure 2: Constitutively expressed genes are characterized by low entropy values**

342     A) The relative frequency of entropy values for a list of housekeeping genes is shown as a kernel

343     density estimation (KDE) plot. The rug plot, black lines on the bottom in the KDE plot represents

344     the individual data points that create the estimation. The y-axis is the probability density, which

345     is the probability for each unit (gene) on the x-axis. The total area below the KDE curve

346     integrates to one.

347     B) The heatmap shows the expression value for housekeeping genes across all conditions

348     analyzed. The expression level for each gene is plotted as the $\log_2$ transformed transcript per

349     million (TPM) value. Genes (rows) are plotted in ranked order based on the entropy value from

350     low (top) to high (bottom). The scale on the left indicates entropy values for each gene. Each

351     condition (column) has been assigned a category: Metabolism (gold), Development (green), or

352     Light Response (blue). The categories are represented at the top of the heatmap in the three

353     different colors.

354    C) The relative frequency of entropy values for a list of genes related to transcription and

355    translation is shown as a kernel density estimation (KDE) plot. The rug plot, black lines on the

356    bottom in the KDE plot represents the individual data points that create the estimation. The y-

357    axis is the probability density, which is the probability for each unit (gene) on the x-axis.

358    D) Heatmap of $\log_2$ transformed TPM values from all transcription and translation related genes

359    (rows) ranked by entropy (low to high). Entropy values are depicted by the brown to green

360    heatmap on the left, where brown is low (top) and green is high (bottom). Each condition

361    (column) has been assigned a category: Metabolism (gold), Development (green), or Light

362    Response (blue). The categories are represented at the top of the heatmap in the three

363    different colors.

364

365    **Figure 3: Validating entropy values with previously published light induced genes and genes**

366    **induced during sexual development**

367    A) The relative frequency of entropy values for a list of light induced genes is shown as a kernel

368    density estimation (KDE) plot. The rug plot, black lines on the bottom in the KDE plot represents

369    the individual data points that create the estimation. The y-axis is the probability density, which

370    is the probability for each unit (gene) on the x-axis. The total area below the KDE curve

371    integrates to one.

372    B) The heatmap shows the expression value for light induced genes across all conditions

373    analyzed. The expression level for each gene is plotted as the $\log_2$ transformed TPM value.

374    Genes (rows) are plotted in ranked order based on the entropy value from low (top) to high

375    (bottom). The scale on the left indicates entropy values for each gene. Each condition (column)

376      has been assigned a category: Metabolism (gold), Development (green), or Light Response

377      (blue). The categories are represented at the top of the heatmap in the three different colors.

378      C) The relative frequency of entropy values for a list of sexual development genes genes is

379      shown as a kernel density estimation (KDE) plot. The rug plot, black lines on the bottom in the

380      KDE plot represents the individual data points that create the estimation. The y-axis is the

381      probability density, which is the probability for each unit (gene) on the x-axis. The total area

382      below the KDE curve integrates to one.

383      D) The heatmap shows the expression value for sexual development genes across all conditions

384      analyzed. The expression level for each gene is plotted as the $\log_2$ transformed TPM value.

385      Genes (rows) are plotted in ranked order based on the entropy value from low (top) to high

386      (bottom). The scale on the left indicates entropy values for each gene. Each condition (column)

387      has been assigned a category: Metabolism (gold), Development (green), or Light Response

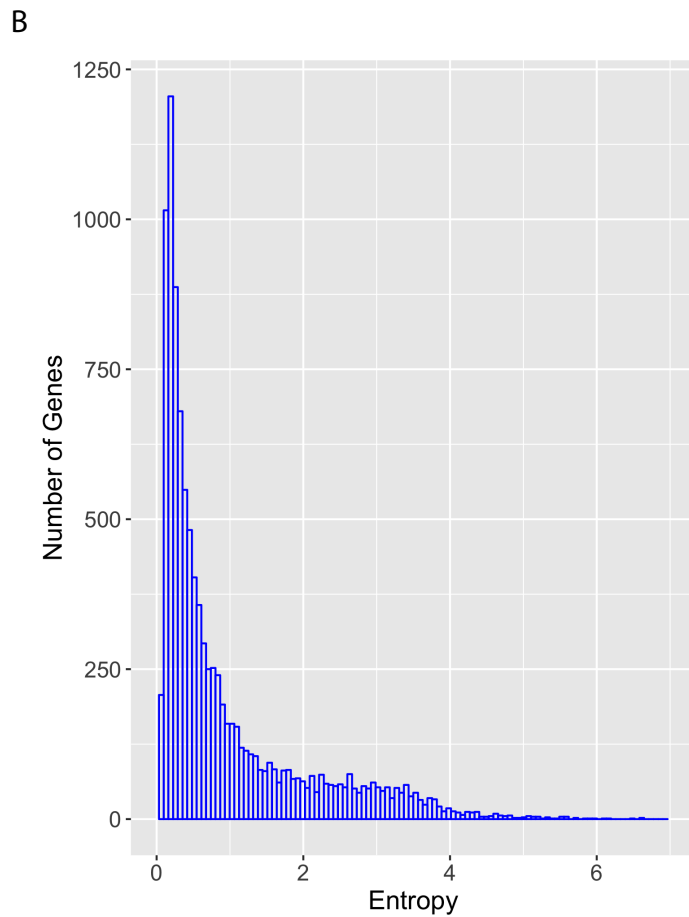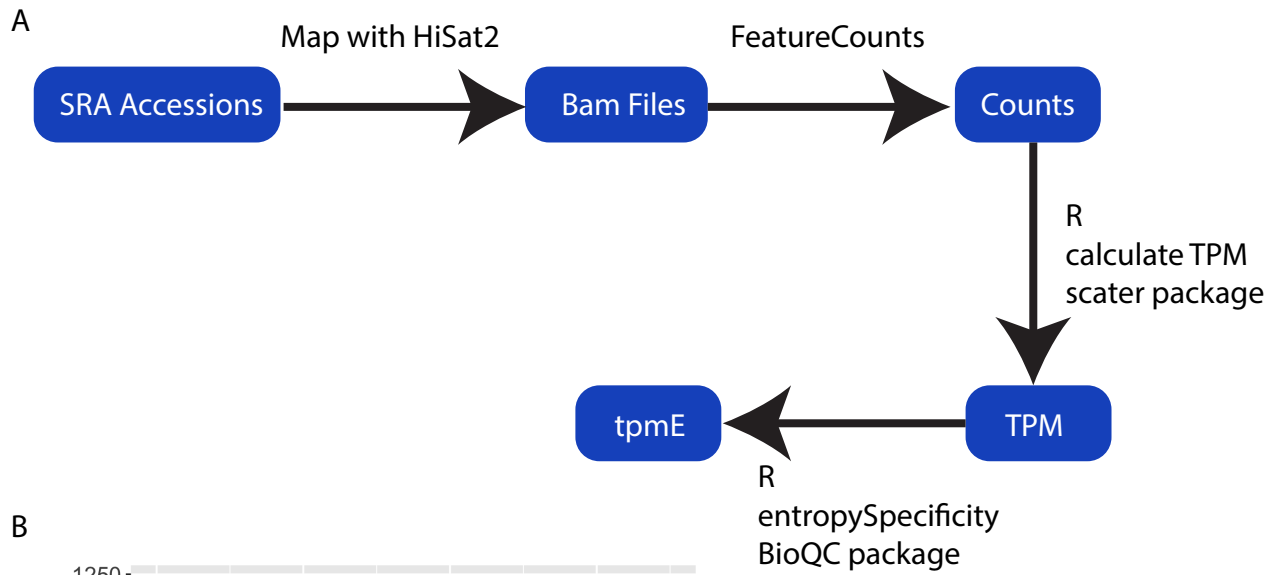388      (blue). The categories are represented at the top of the heatmap in the three different colors.

389

390      **Figure 4: Log$_2$ TPM values for highest and lowest ranked genes**

391      A) The heatmap shows the expression values for the 100 genes with the highest entropy values.

392      The expression level for each gene is plotted as the $\log_2$ transformed TPM value. Each row

393      represents a gene. Gene names are listed on the right side of the heatmap. Each condition

394      (column) has been assigned a category: Metabolism (gold), Development (green), or Light

395      Response (blue). The categories are represented at the top of the heatmap in the three
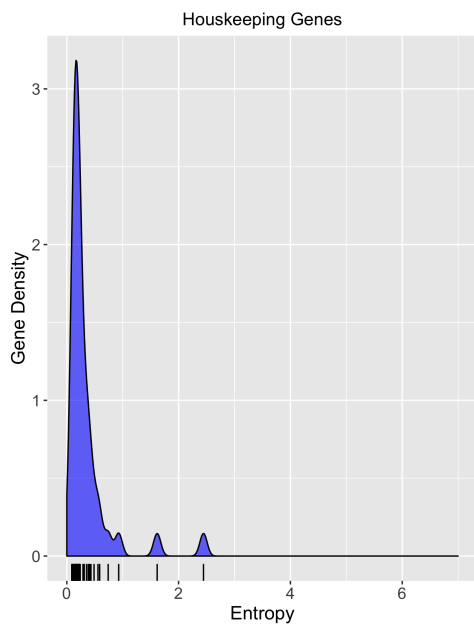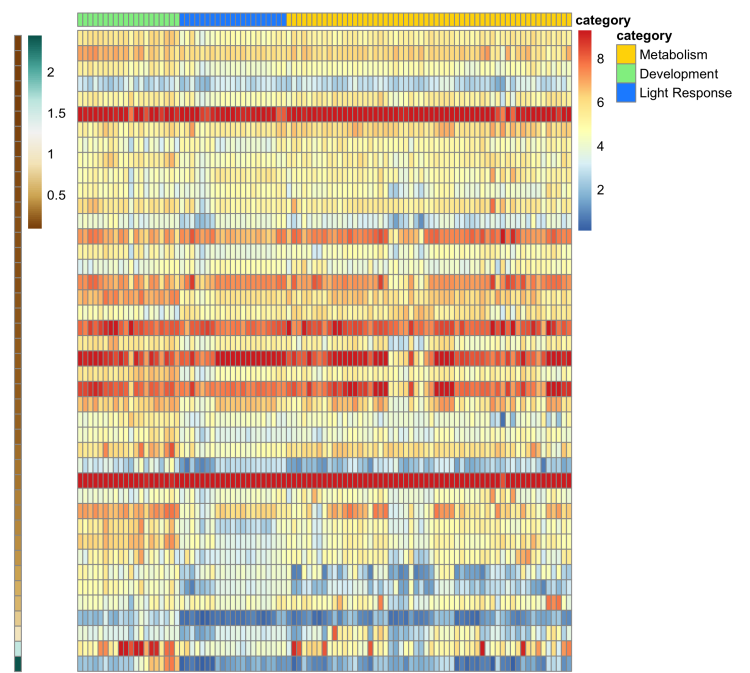
396      different colors.

397    B) The heatmap shows the expression values for the 100 genes with the lowest entropy values.

398    The expression level for each gene is plotted as the $\log_2$ transformed TPM value. Each row

399    represents a gene. Gene names are listed on the right side of the heatmap. Each condition

400    (column) has been assigned a category: Metabolism (gold), Development (green), or Light

401    Response (blue). The categories are represented at the top of the heatmap in the three

402    different colors.

403
404
405

A



B



C



D

A



B